

Impact of Non-optimal Checkpoint Intervals on Overall Application Efficiency in Cluster Computing: *A Simulation-Based Study*

William M. Jones, PhD
Computer Science Department

Coastal Carolina University
Myrtle Beach, South Carolina

<http://www.parl.clemson.edu/beosim>

The Issue

- Application efficiency
 - Throughput
 - Turnaround times
 - Heavy loads
- Optimal checkpoint intervals
- AMTTI estimation
 - Running estimate
- Sensitivity to error
 - Non-optimal checkpoint interval
- Simulation-based study
 - Discrete event-driven
 - LANL's Pink Cluster

Parameter Estimation



Checkpoint Interval



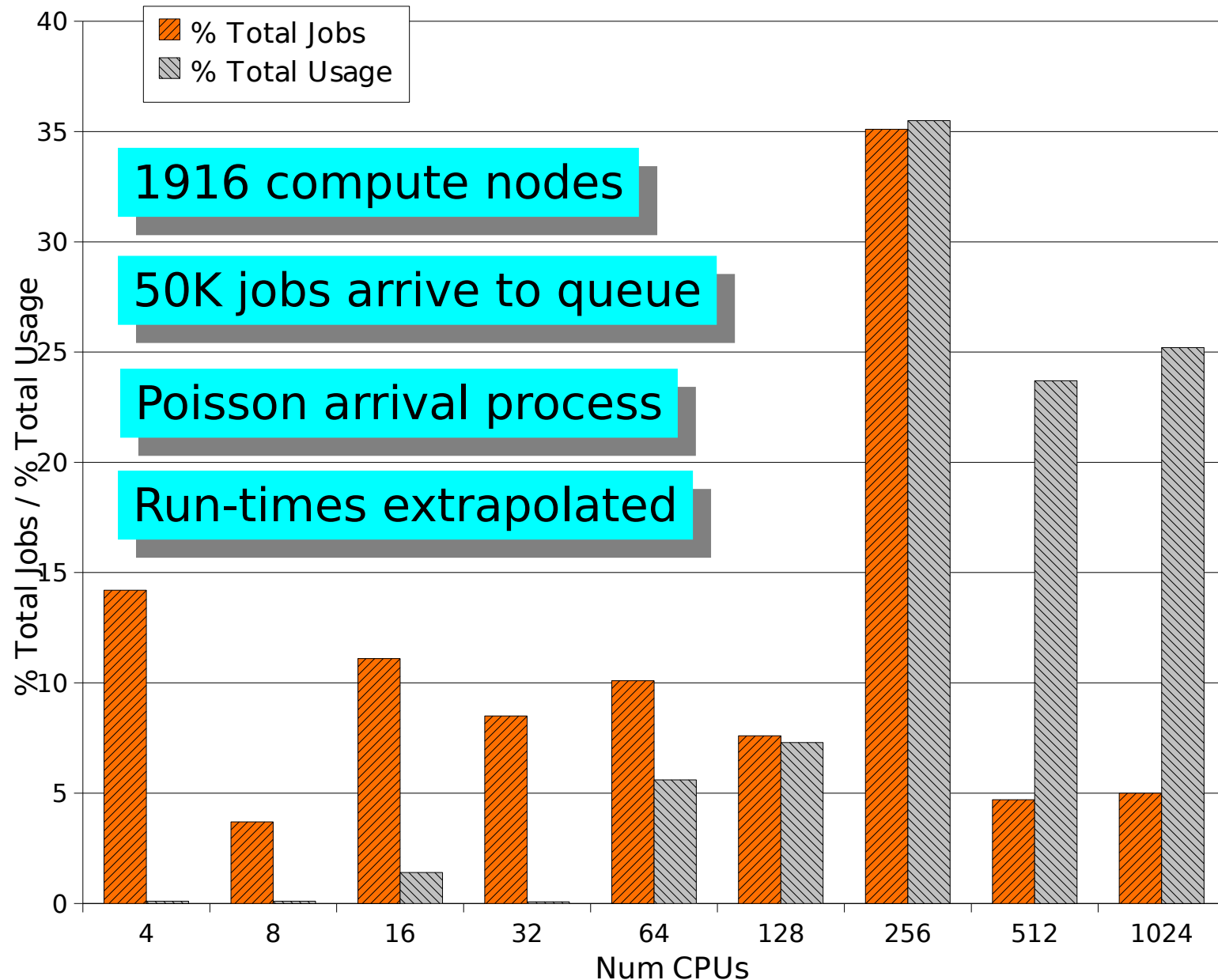
Application Efficiency



Throughput, TAT, etc

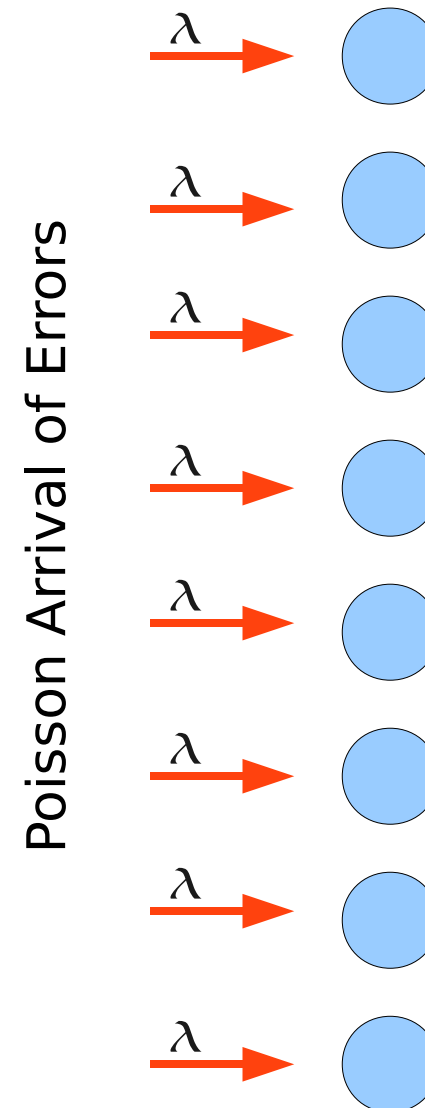
LANL's PINK Cluster Workload

PINK (01/2007 to 11/2007) Job Stats

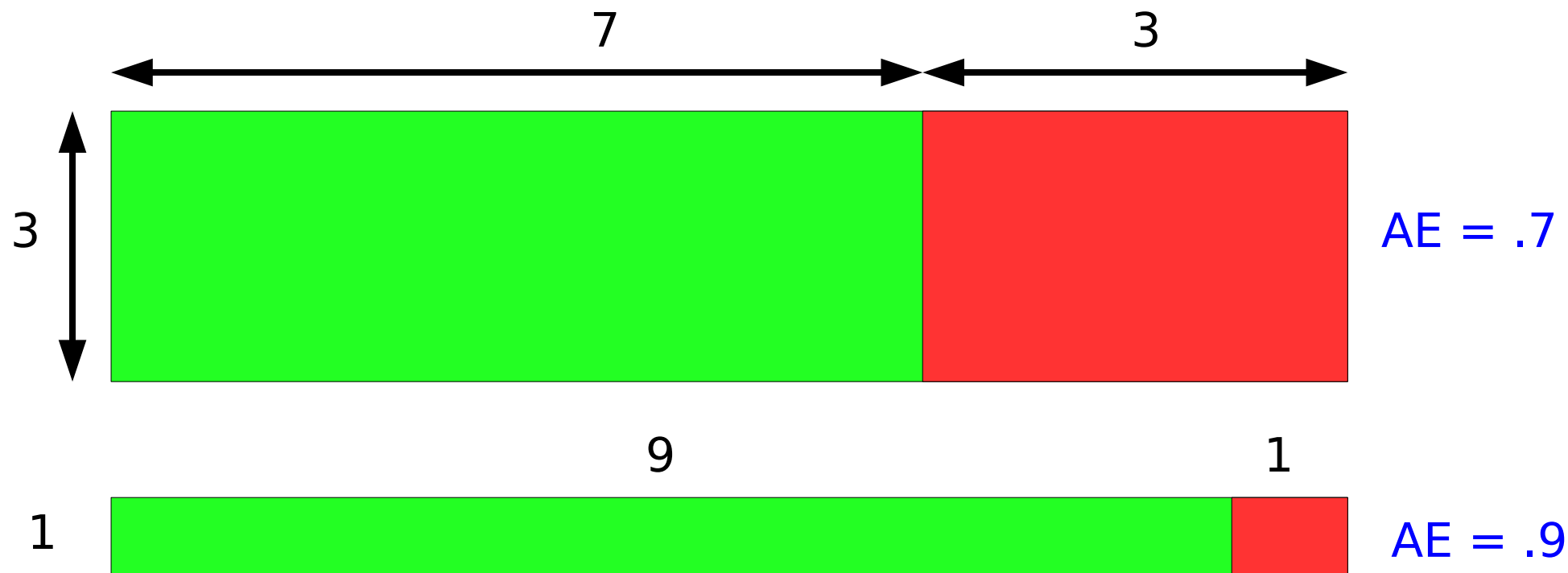


PINK Failure Model (Assumptions)

- Independent
- MTBF / node
- Single failure mode
- One failure
 - Halts at most 1 job
- Zero repair time
 - Job added to HOQ
- Checkpoint / Restart
 - 10 min dump time
 - 10 min restart time
- Simplistic
 - First pass



Checkpointing Metrics: Average Application Efficiency



$$Avg. AE = \frac{\sum_{i=0}^{N-1} Ts_i n_i}{\sum_{i=0}^{N-1} Tr_i n_i} = .75$$

$$Avg. AE = \frac{\sum_{i=0}^{N-1} \frac{Ts_i}{Tr_i}}{N} = .8$$

Setting Checkpoint Interval For Each Job in Simulation

$$AMTTI_i = \frac{T}{n_i}, T = \frac{1}{\lambda}$$

Assume this

$$Tc_i = \sqrt{2 \delta AMTTI_i}$$

Approximate that

$$Tc_i = Tc_i * Err$$

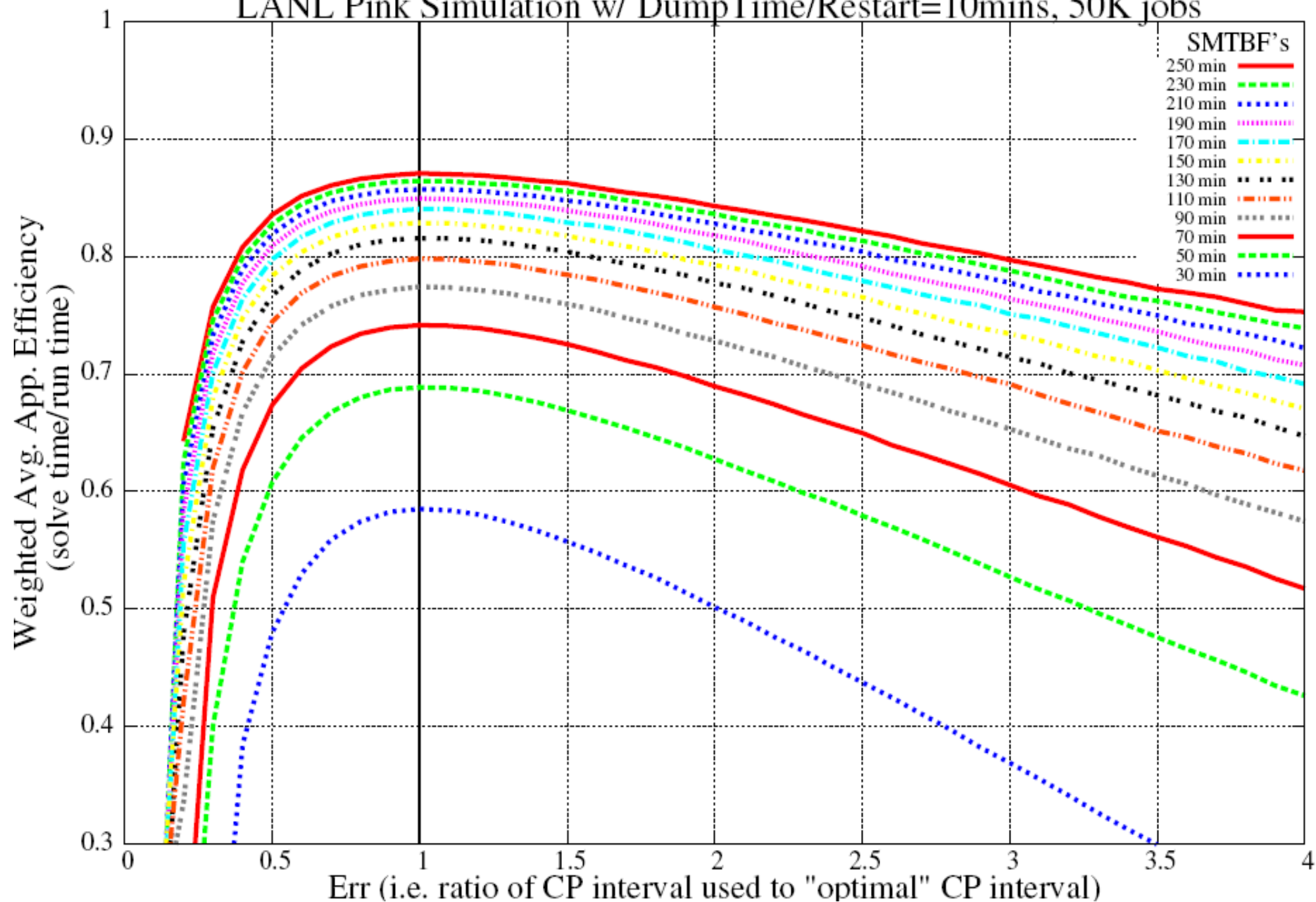
Introduce error here

$$(Tc_i < \delta) \vee (Tc_i > Ts_i) \rightarrow no\ cp'ing$$

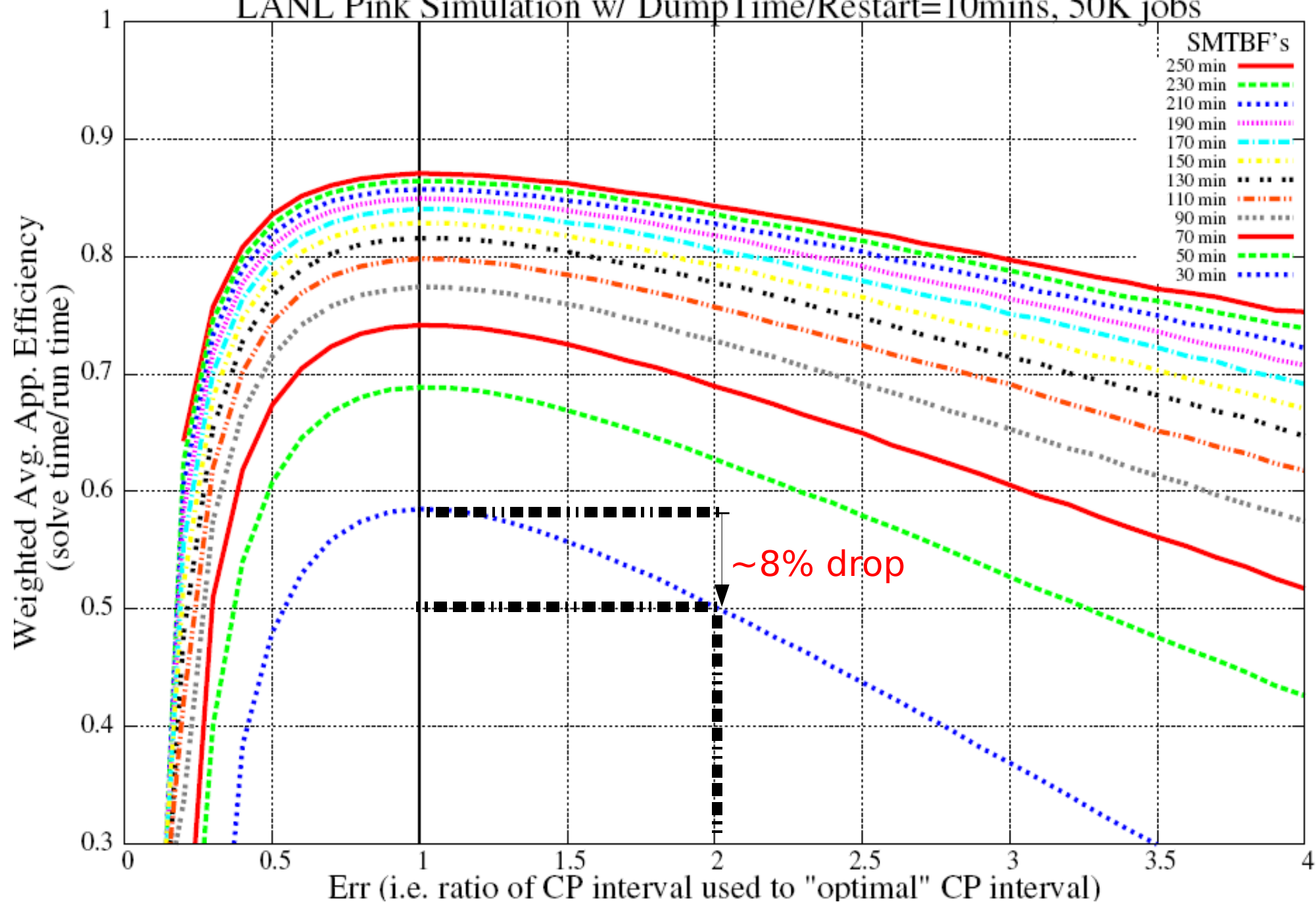
Is Tc logical for job?

Application Efficiency vs Checkpoint Interval vs System MTBF

LANL Pink Simulation w/ DumpTime/Restart=10mins, 50K jobs



Application Efficiency vs Checkpoint Interval vs System MTBF LANL Pink Simulation w/ DumpTime/Restart=10mins, 50K jobs



Back Calculate Loss of Application Efficiency as Function of Error in AMTTI

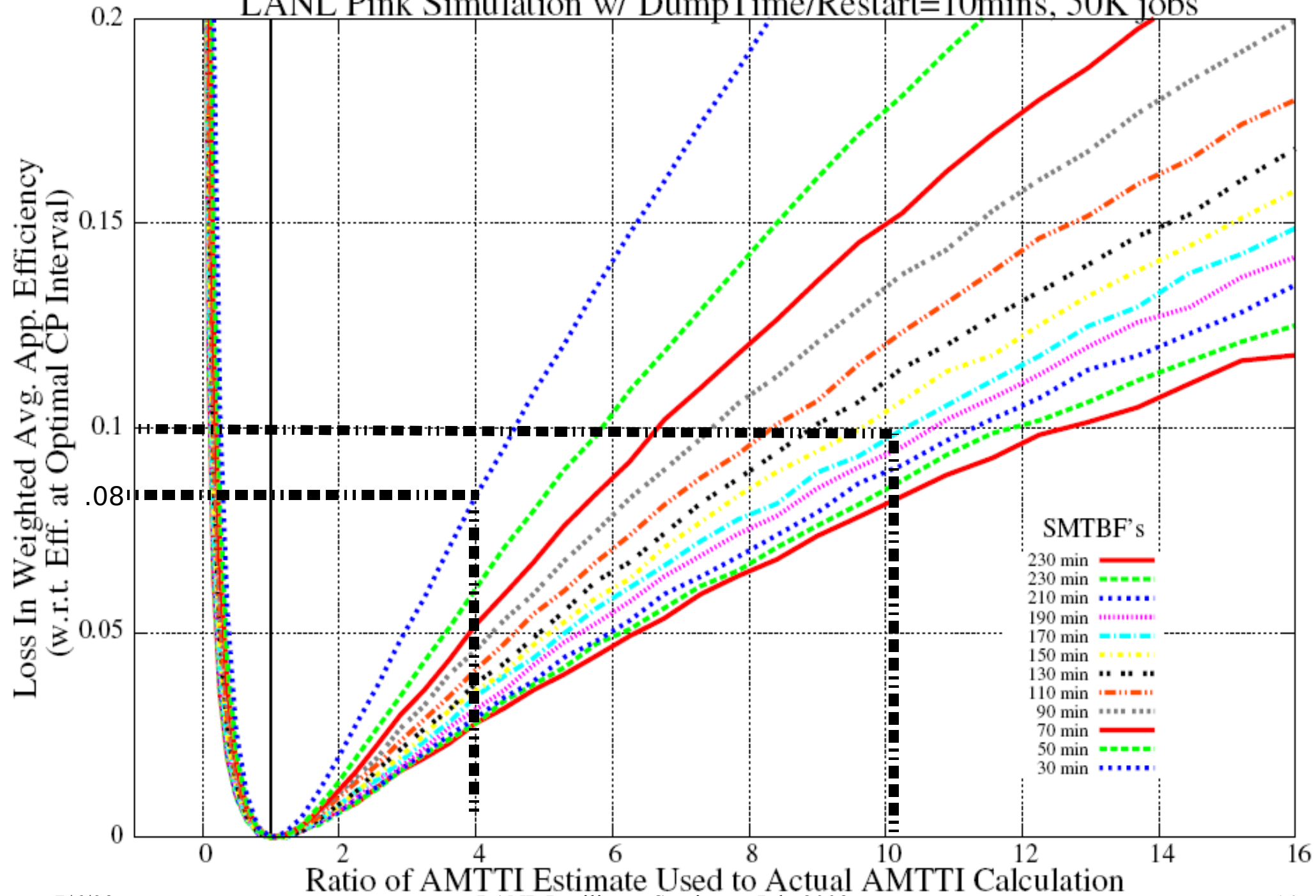
$$Tc_i = \sqrt{2 \delta AMTTI_i}$$

Assume AMTTI is the sole source of error in Tc

A factor of 2 error in Tc would result from a factor of 4 error in AMTTI

Loss In App. Efficiency vs Error In AMTTI Estimate vs System MTBF

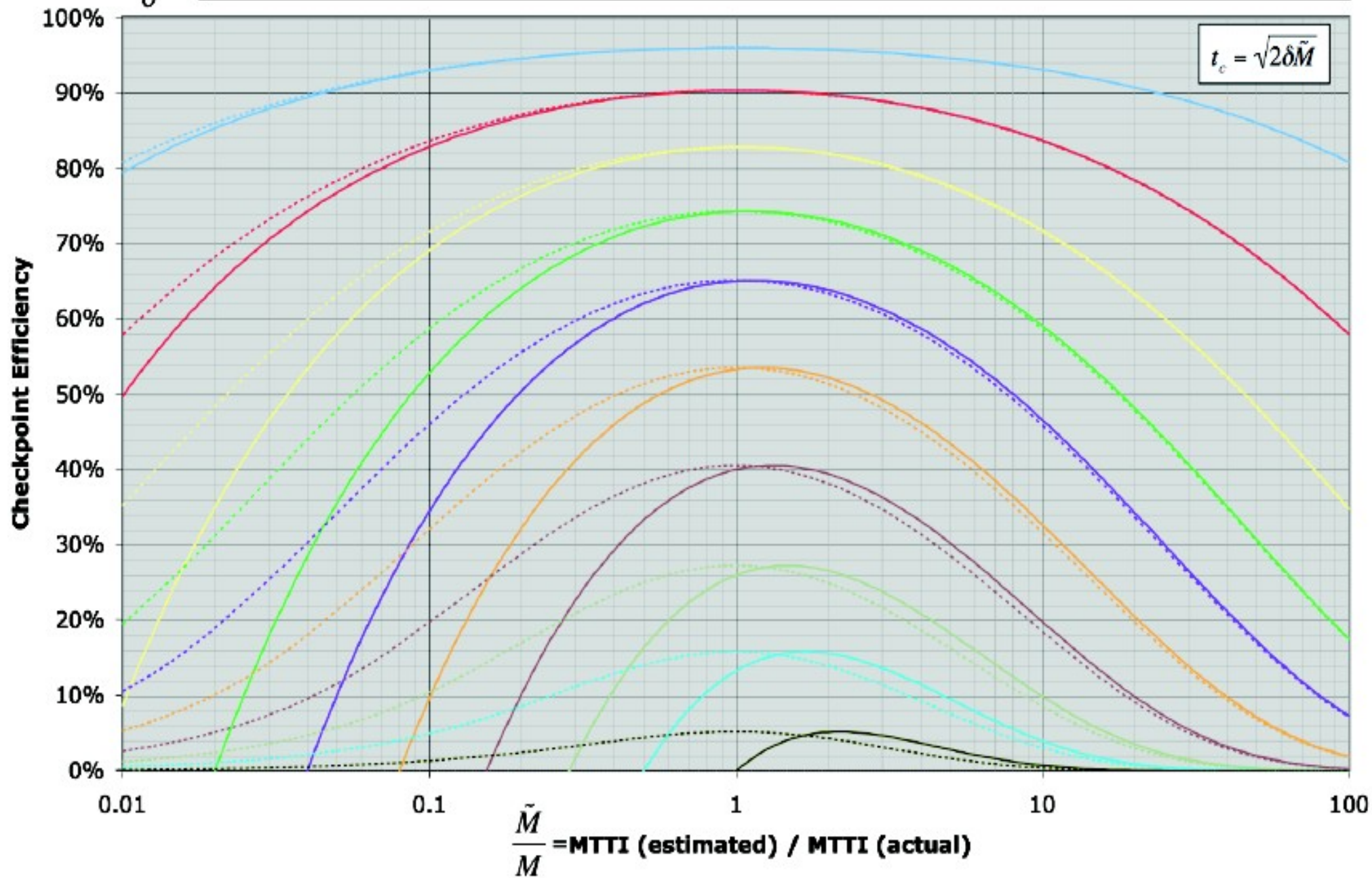
LANL Pink Simulation w/ DumpTime/Restart=10mins, 50K jobs



Sensitivity of an Application to Accuracy of the MTTI Estimate

Solid Lines: Approximate Optimum Checkpoint Interval, Dotted Lines: Exact

$$\frac{M}{\delta} = \begin{matrix} \text{---} 1200 & \text{---} 200 & \text{---} 60 & \text{---} 25 & \text{---} 12.5 & \text{---} 6.25 & \text{---} 3.25 & \text{---} 1.75 & \text{---} 1 & \text{---} 0.5 \end{matrix}$$



Daly, J. T., “Methodology and Metrics for Quantifying Application Throughput”, Proceedings of NECDC ‘06

Daly states “... we do not need to be overly concerned about the precision of our checkpoint restart interval approximation [...] or our measurement of its dependent dump time and application MTTI parameters.”

Application Efficiency Still Matters

Heavy System Load



Dramatic Impact on Queue Time



Turnaround Time = Queue Time + Execution Time

Besides

Energy Usage

But what about at larger scales where dumptime \gg AMTTI ?

Future Work

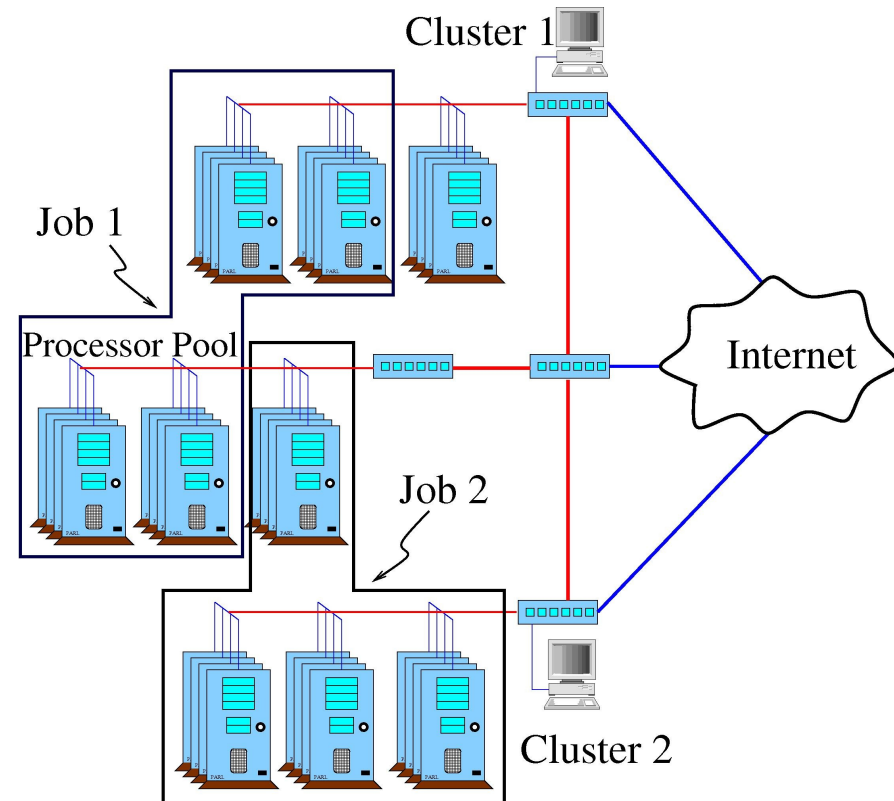
- New cluster (larger systems) runs
 - Realistic failure stats
- Add new failure modes
 - Failures that impact more than one job
 - AMTTI \ll SMTBF
- Add AMTTI on-line parameter estimator
 - Initialize to a guess
 - As AMTTI is refined, T_c becomes closer to optimal
- Map jobs with resilience in mind
 - Schedule contiguous across switch, etc

PART 2 of Presentation

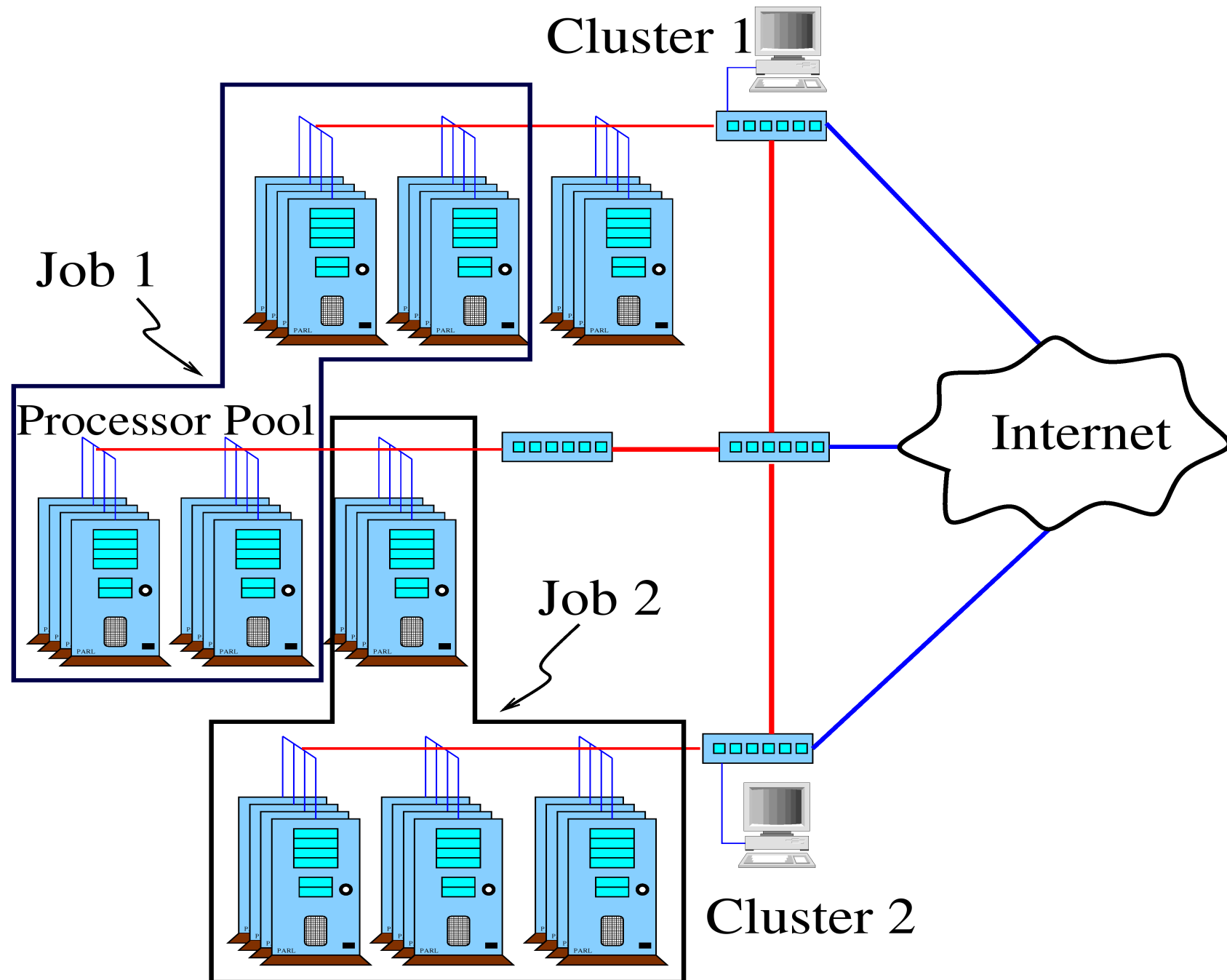
Checkpoint / Migration Job Scheduling Next

Parallel Job Scheduling

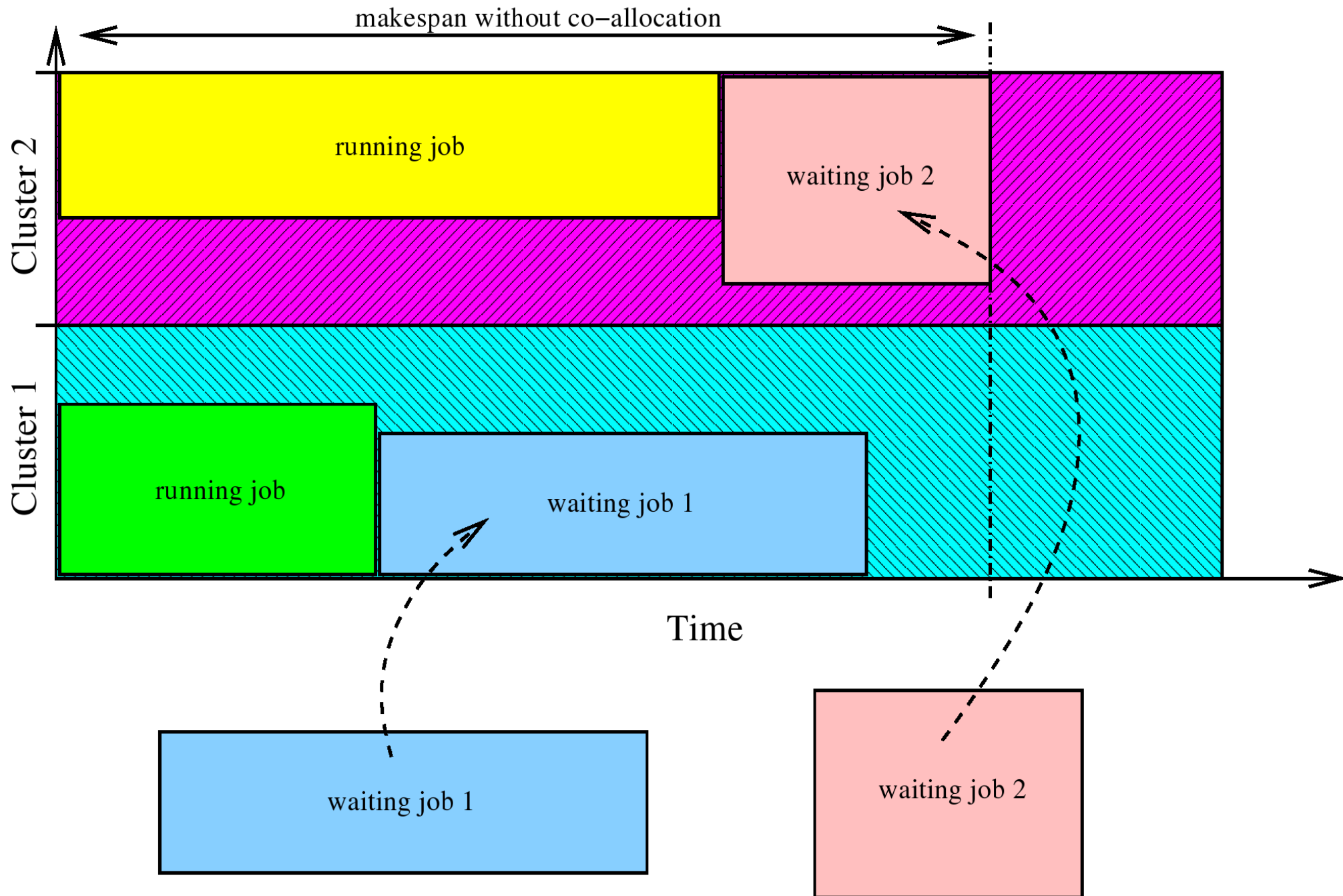
- Local execution
- Job migration
- **Job co-allocation**
 - Map across boundaries
 - Sharing resources
 - Network BW contention
- Can help or hurt
- Some example scenarios ...



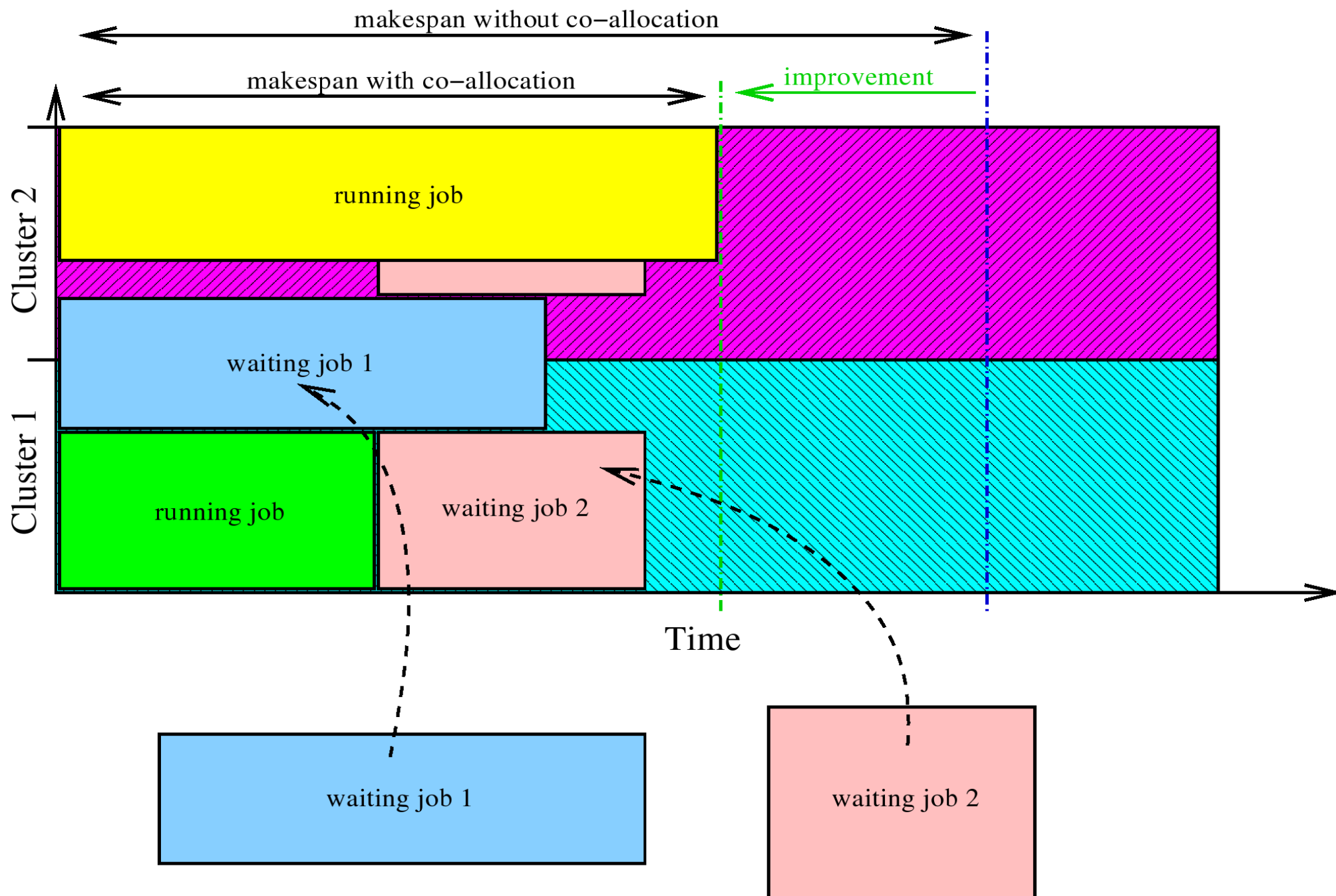
Multi-site Co-allocation



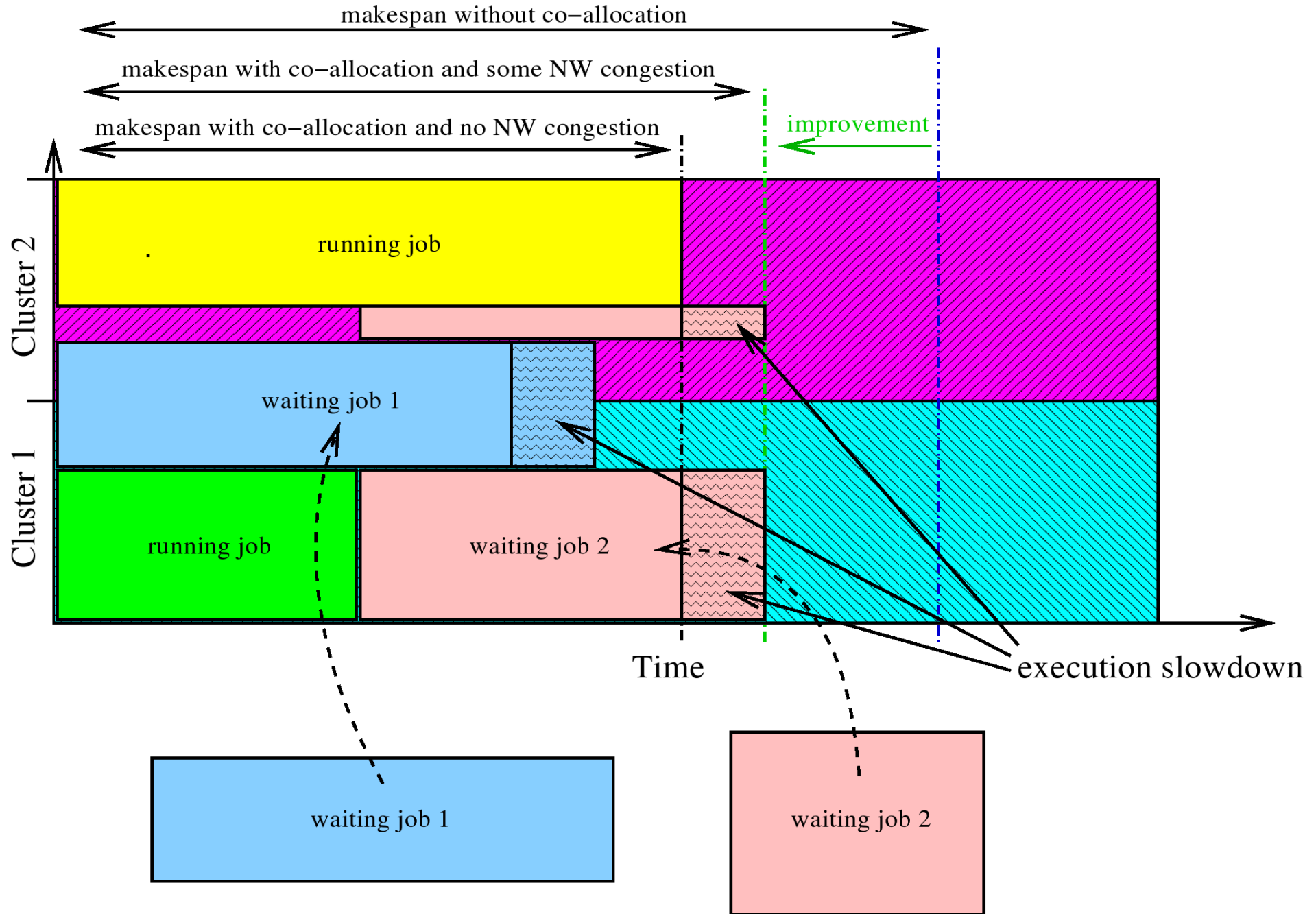
Scheduling w/o co-allocation



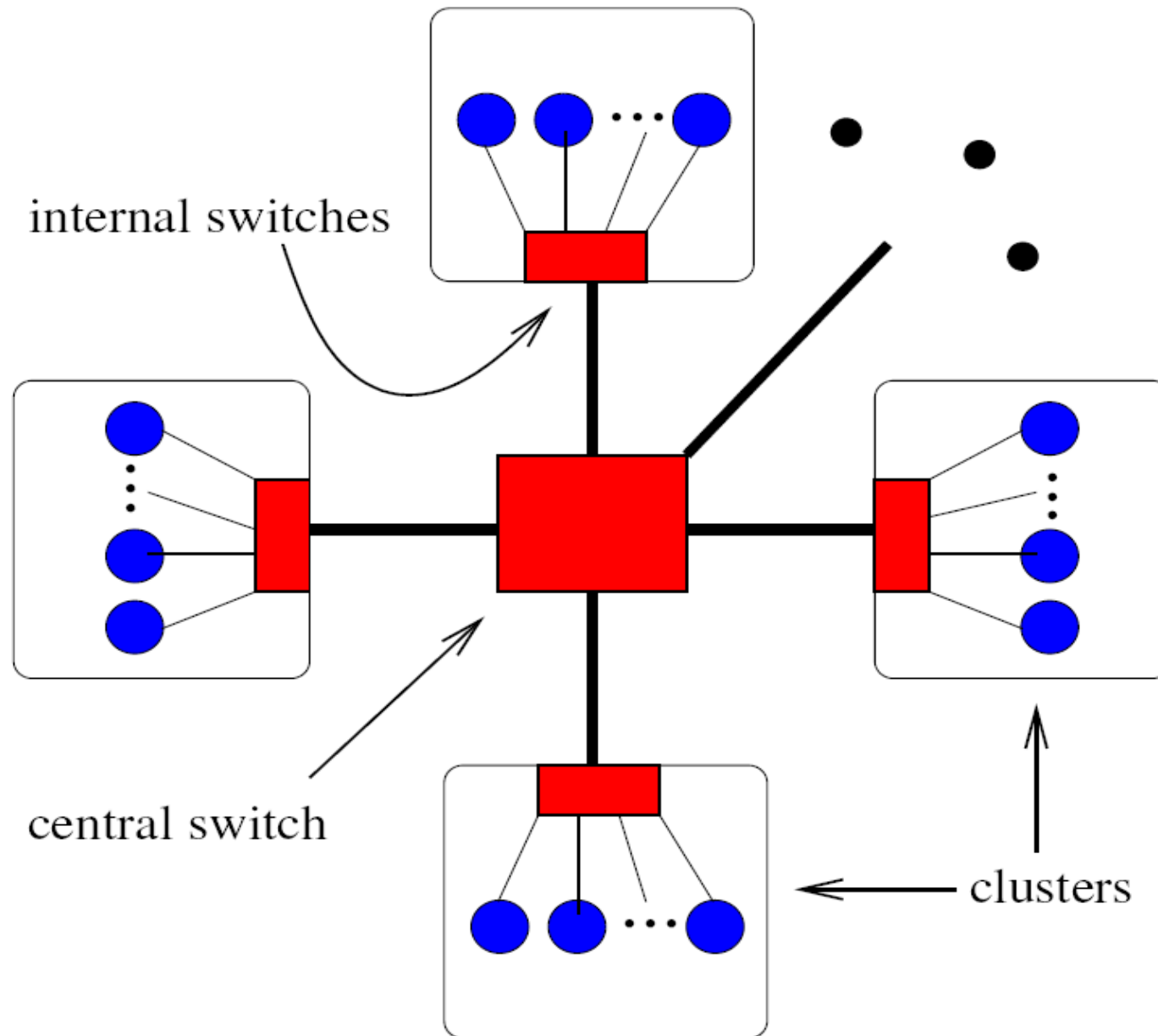
Scheduling w/ co-allocation



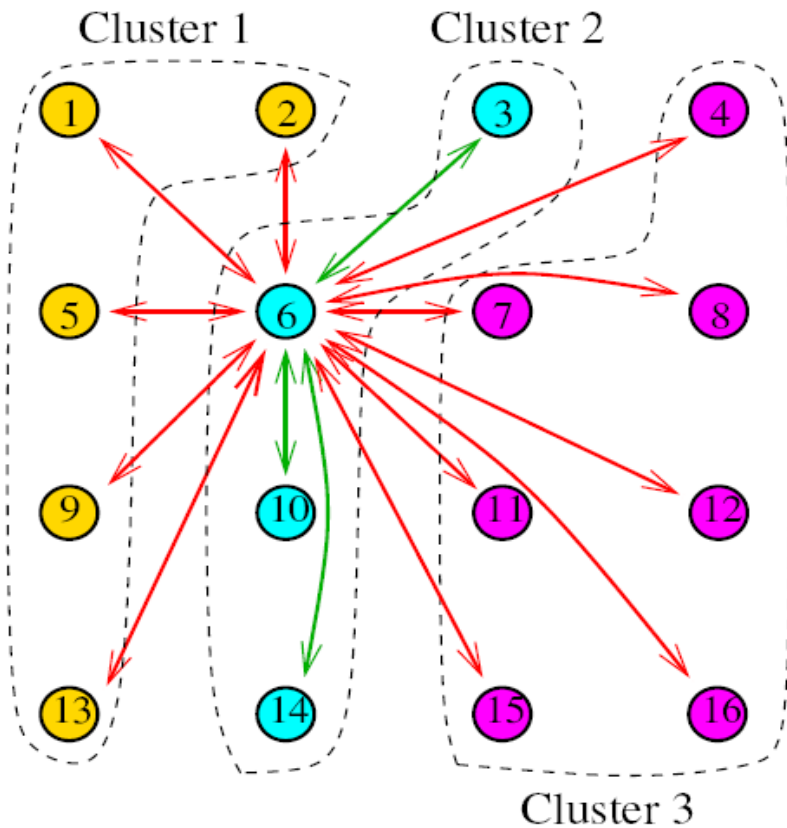
Co-allocation w/ slowdown



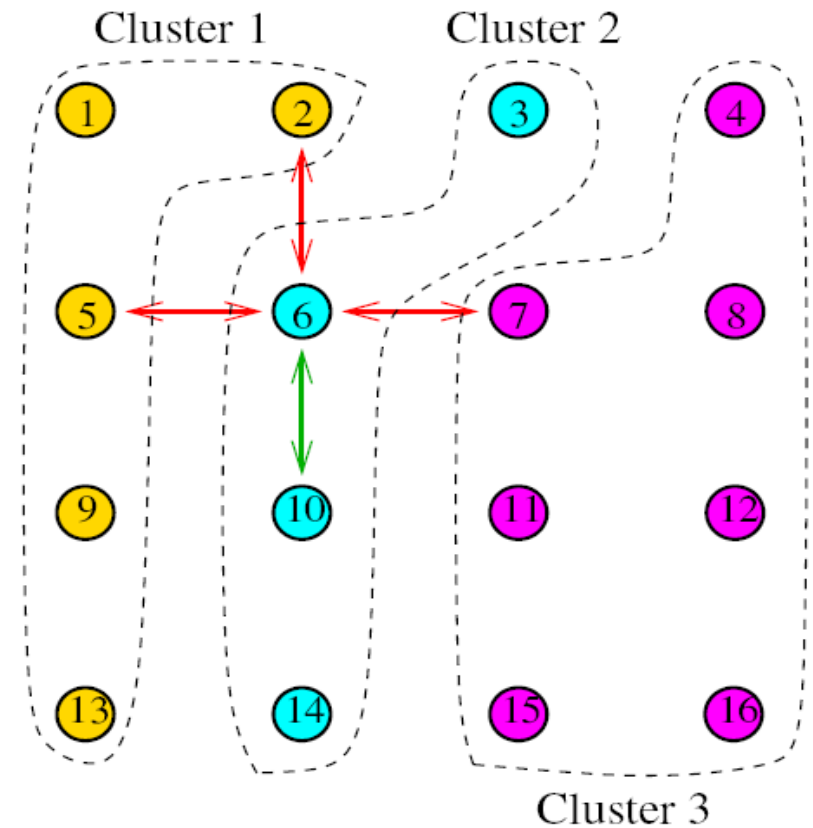
Idealized Multi-cluster Model



IPC Pattern Model



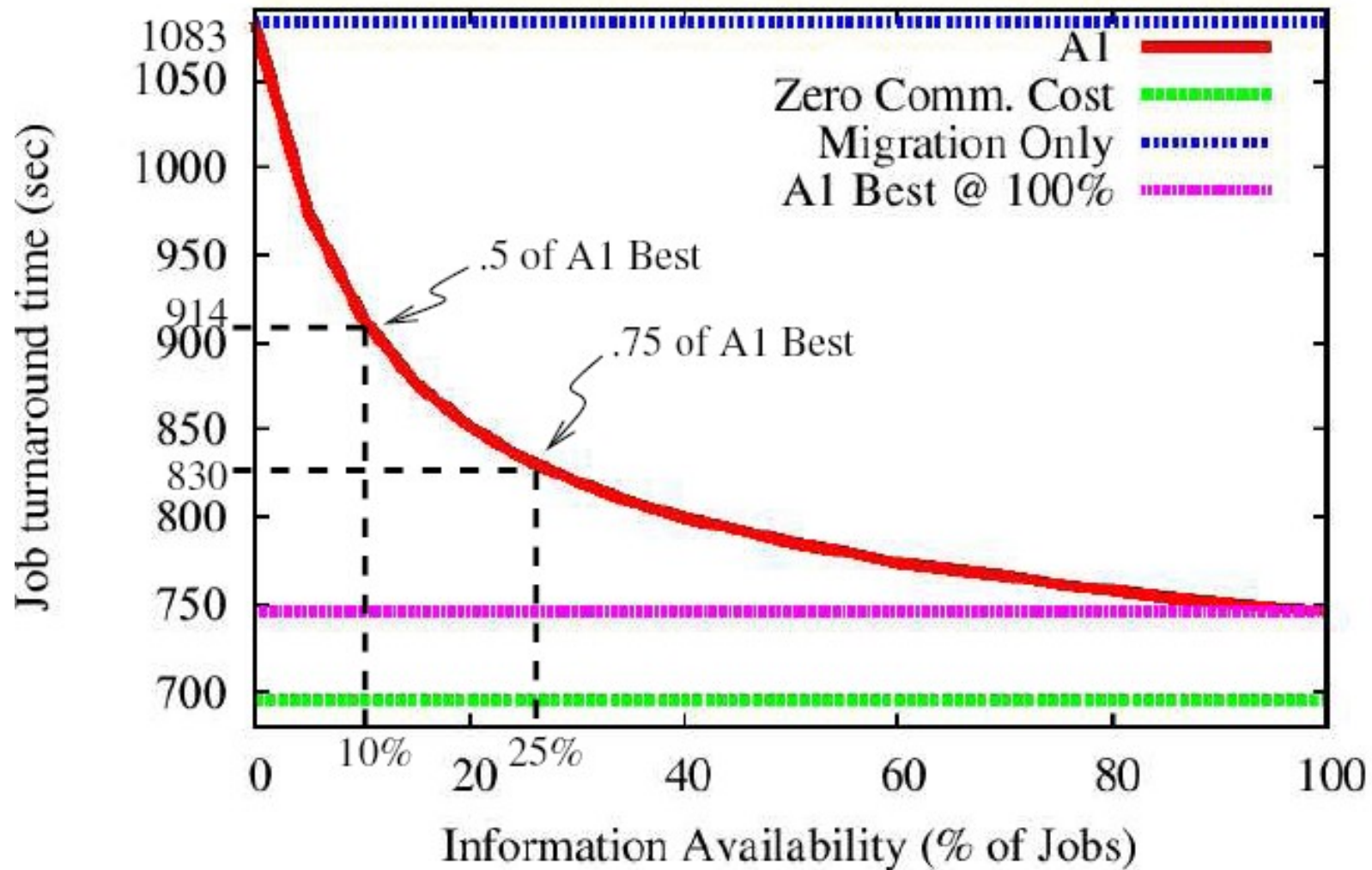
All-to-all personalized



2D

How Much Info Is “Enough”

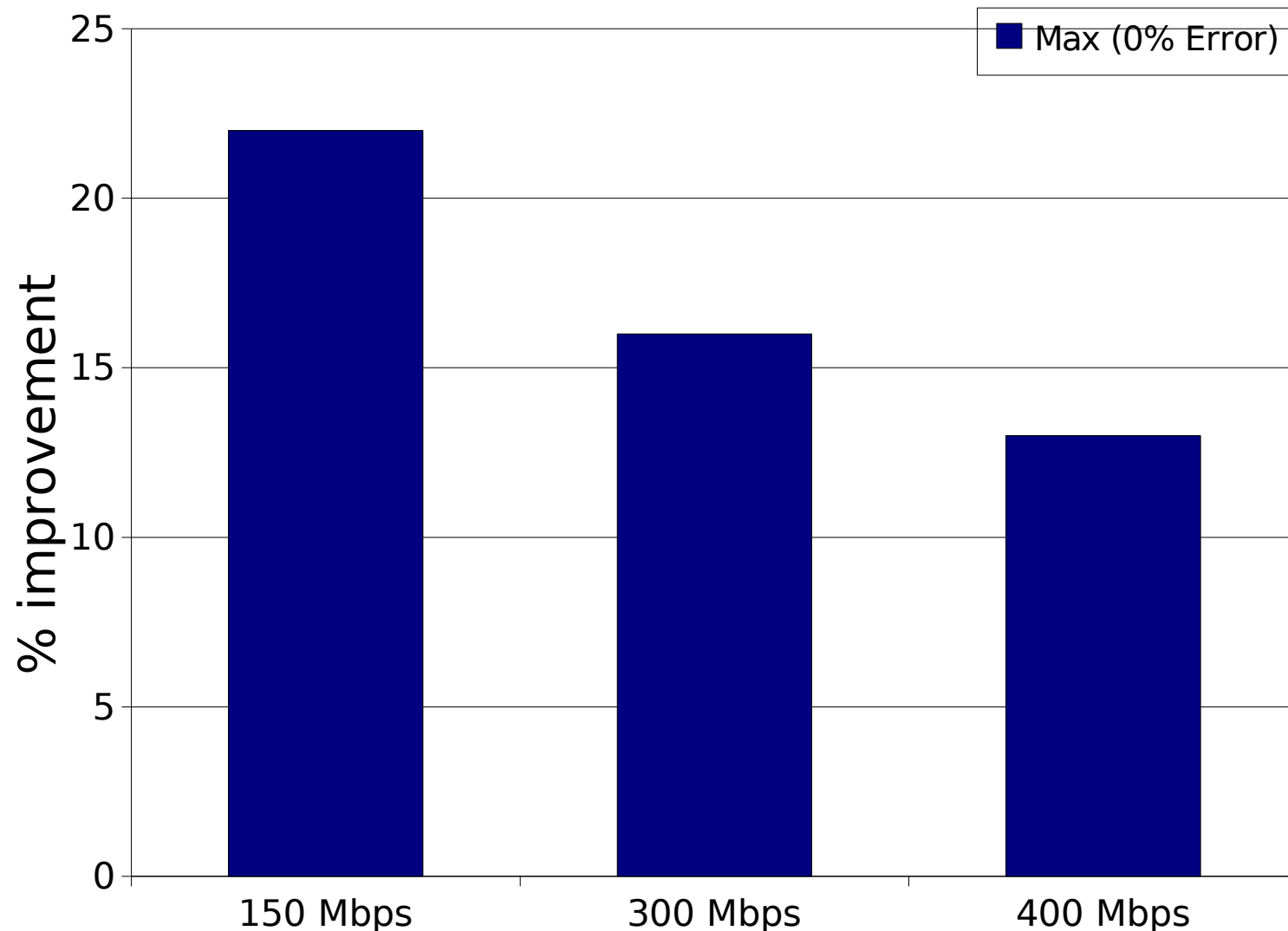
A1 performance w.r.t. Information Availability



Max improvement possible: 36%
A1 Best: 31%

Summary of Previous Results

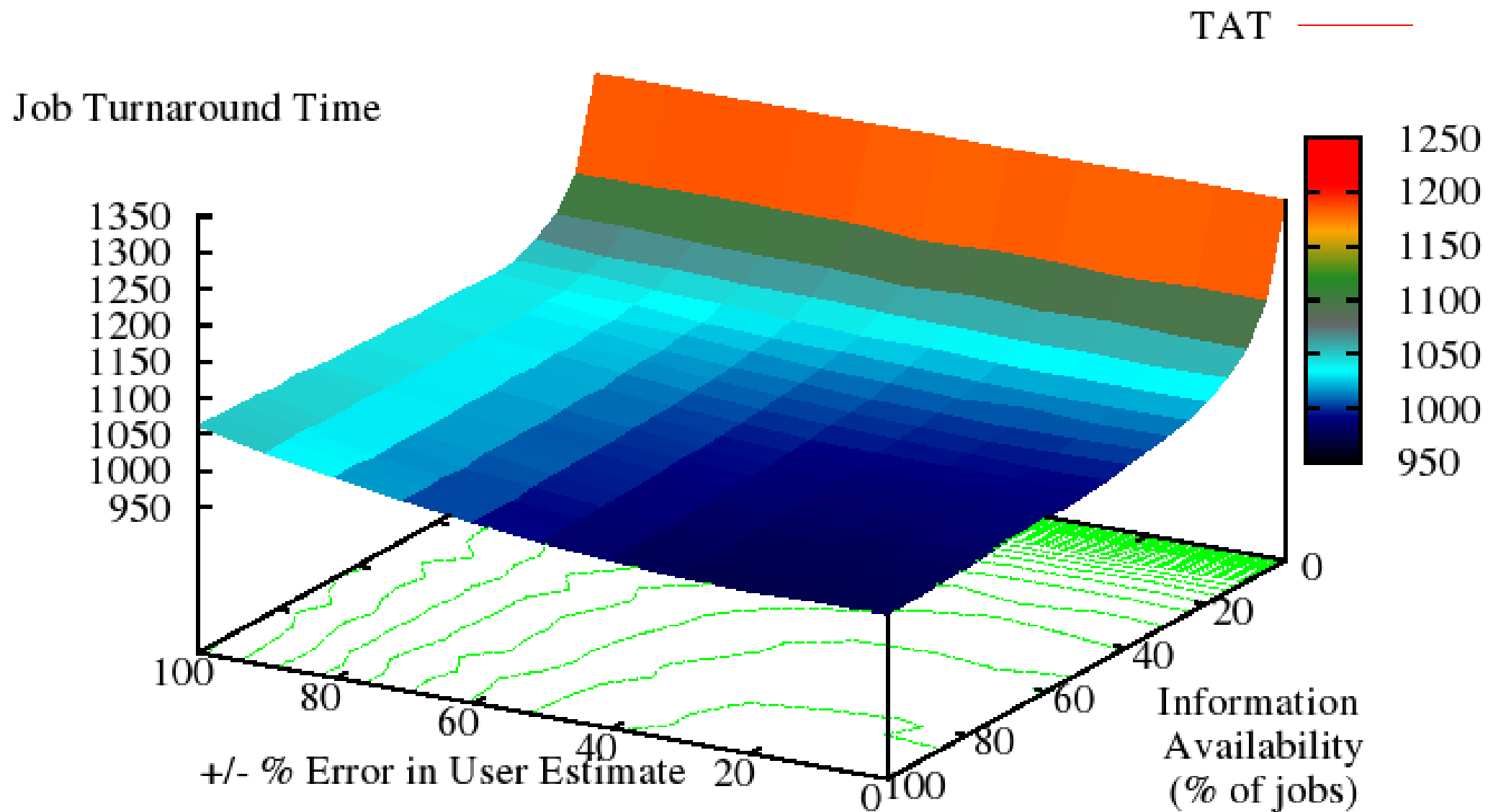
Improvement Over Migration Only w/ No Estimate Inaccuracy



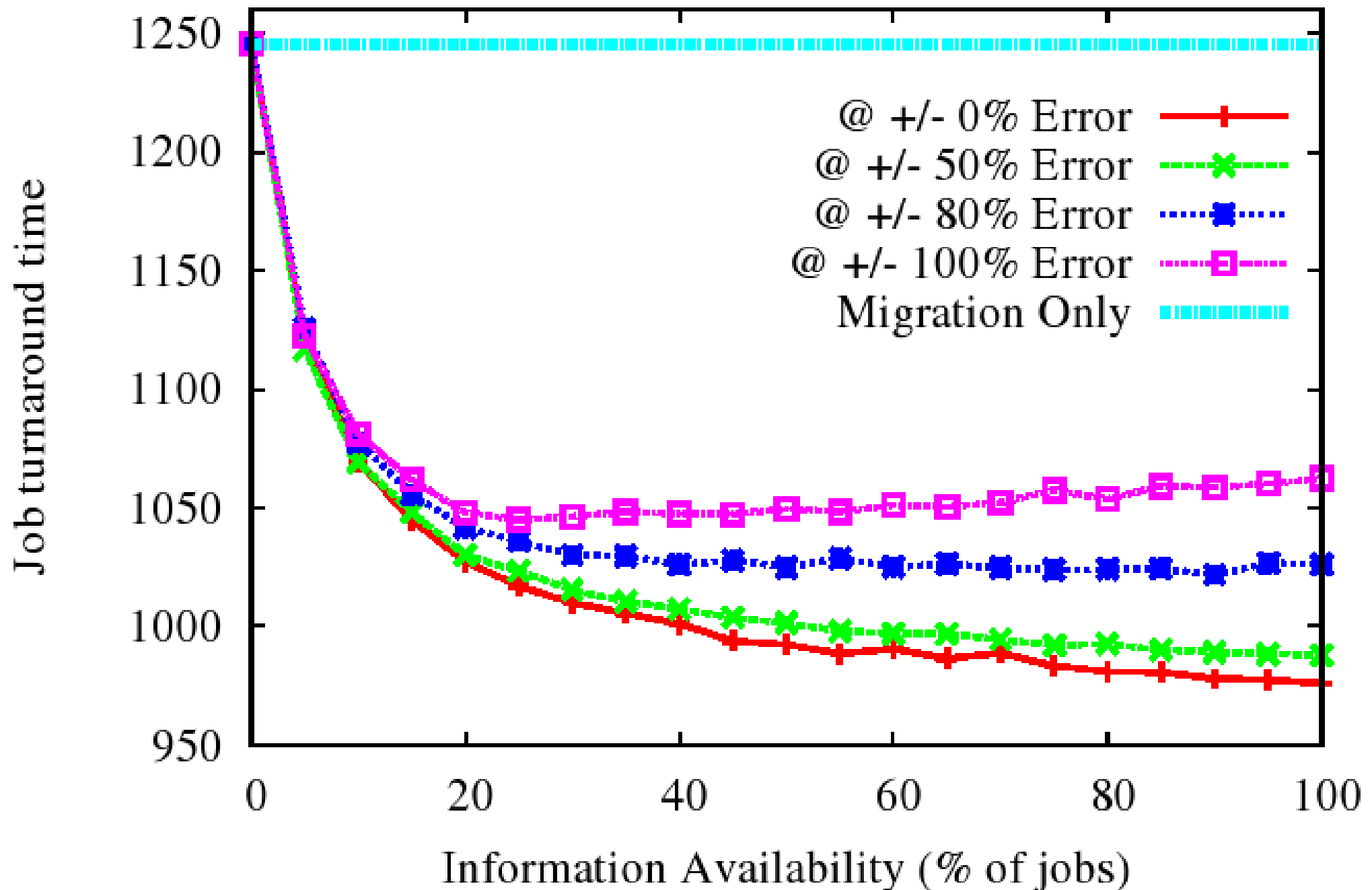


What happens if there is significant inaccuracy in the user-provided bandwidth estimates?

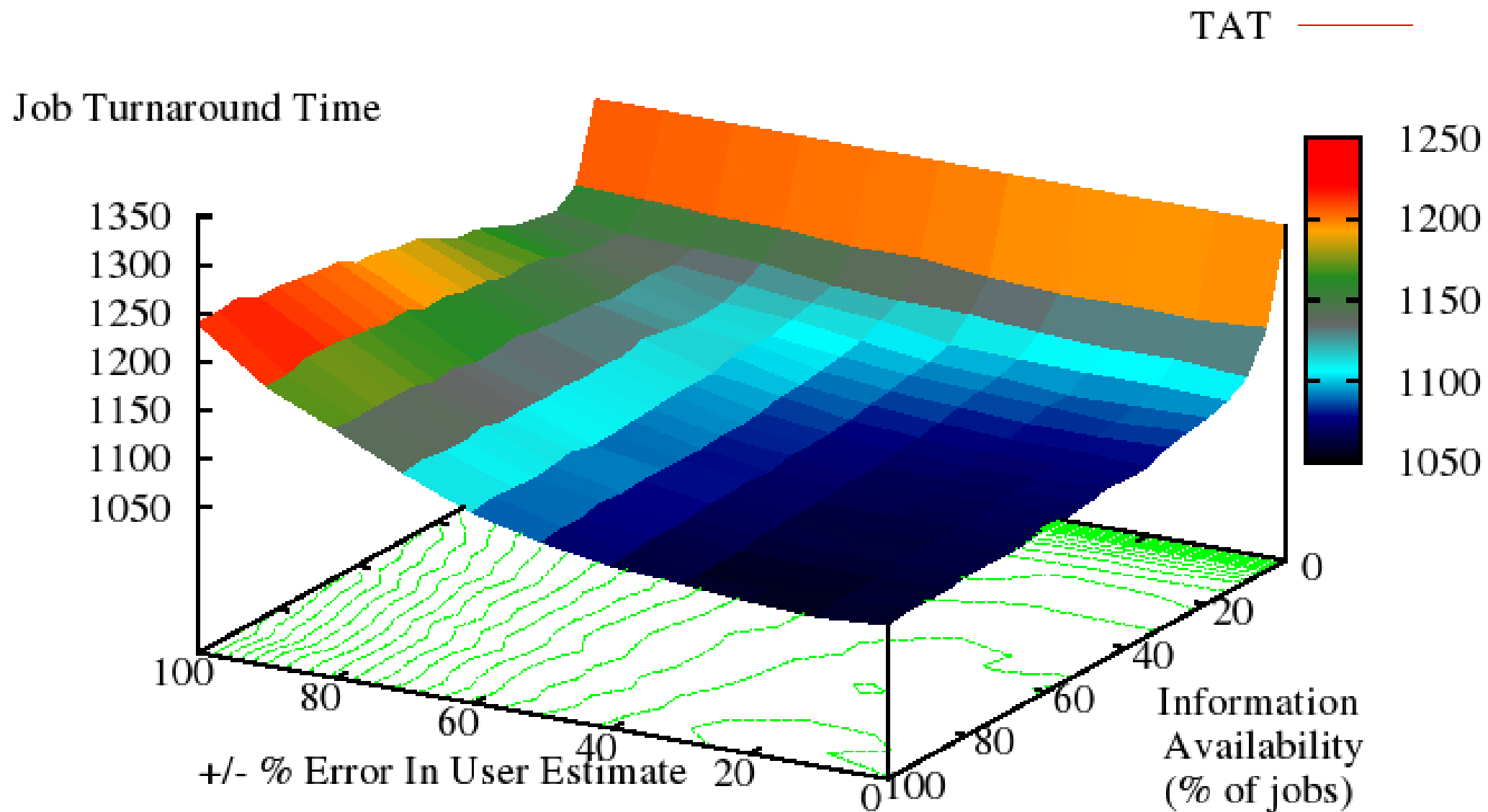
Turnaround Time vs Info vs Error w/ PPBW = 150 Mbps



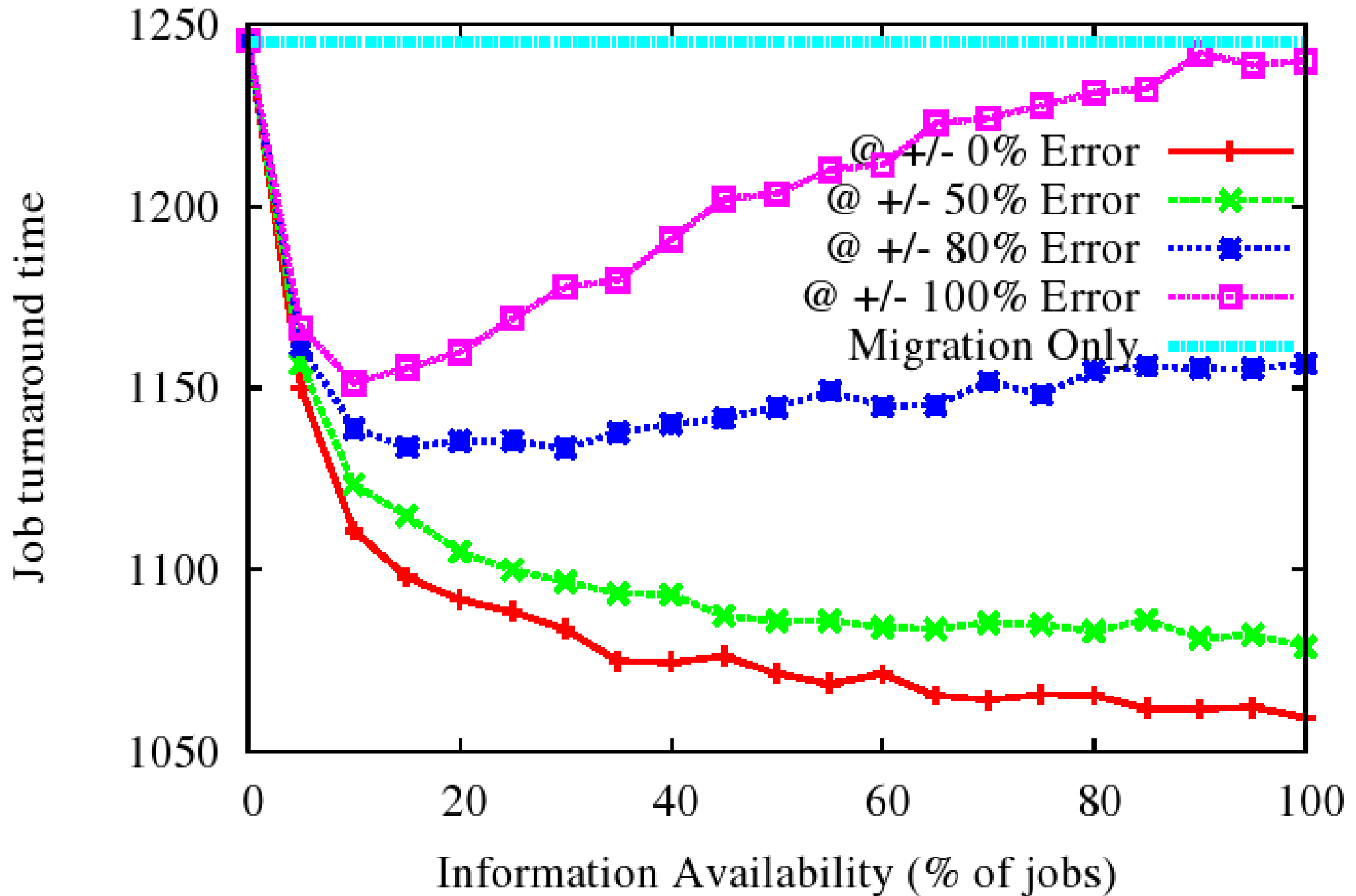
Turnaround Time vs Info w/ PPBW = 150 Mbps



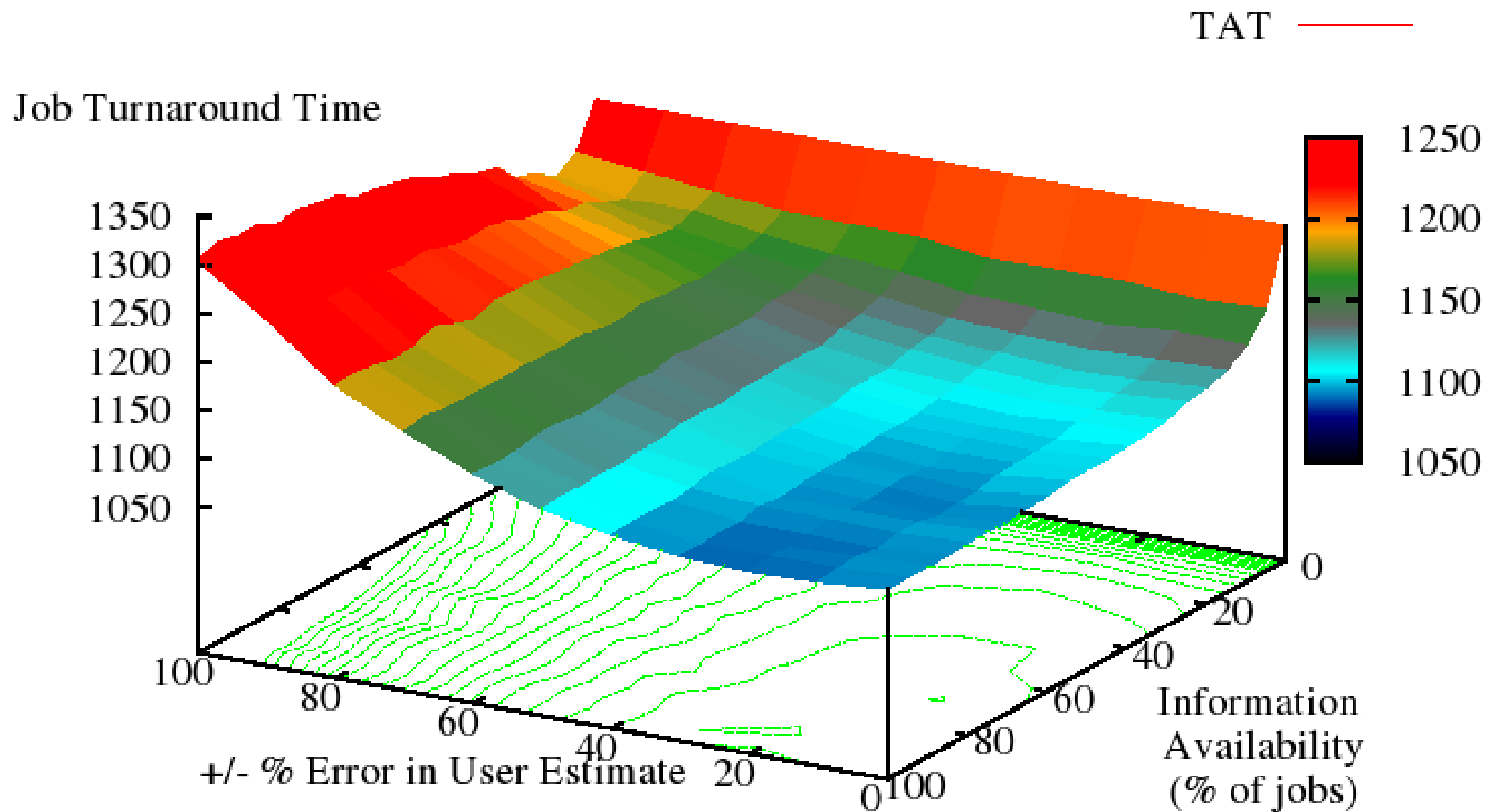
Turnaround Time vs Info vs Error w/ PPBW = 300



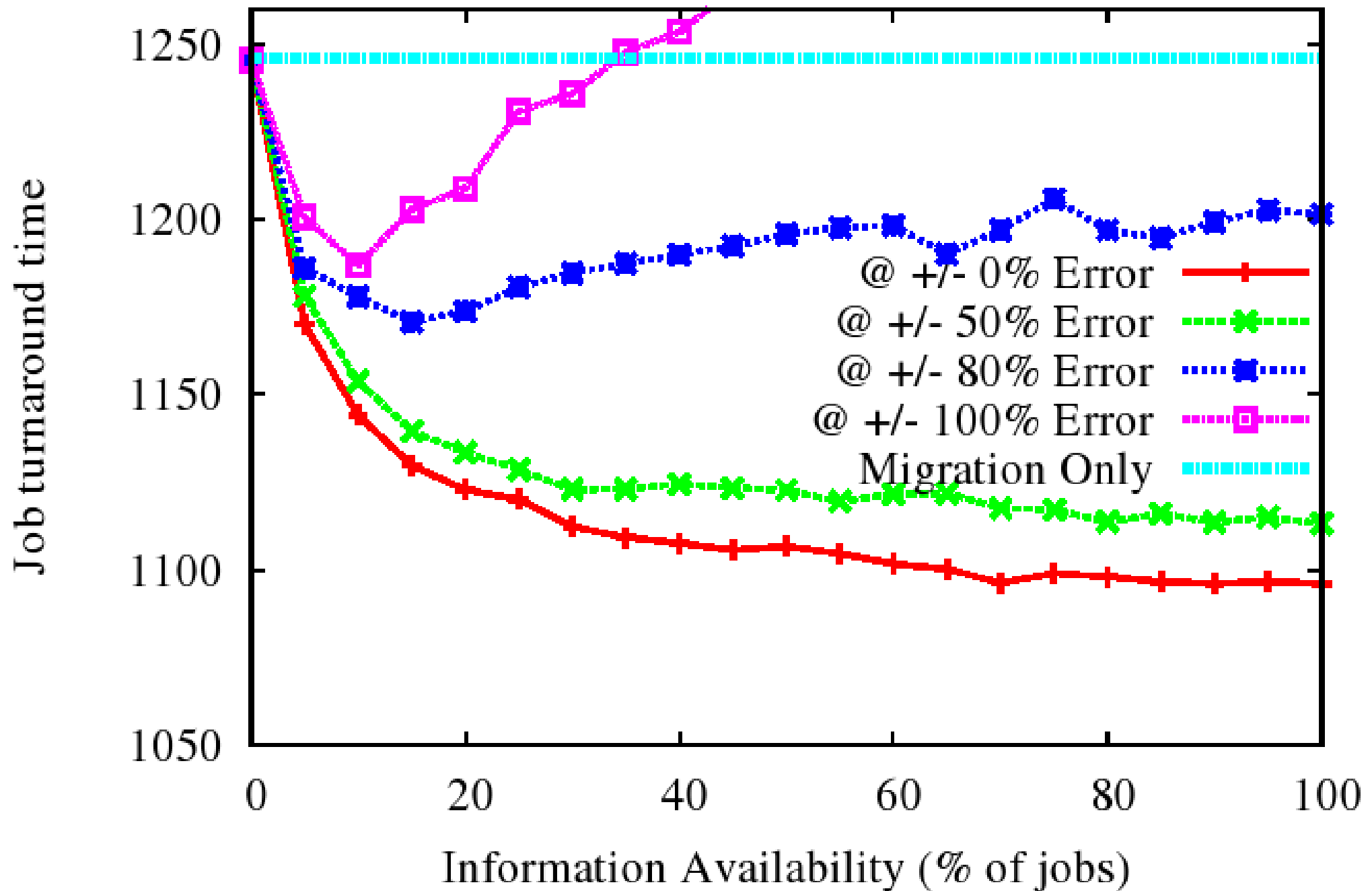
Turnaround Time vs Info w/ PPBW = 300 Mbps



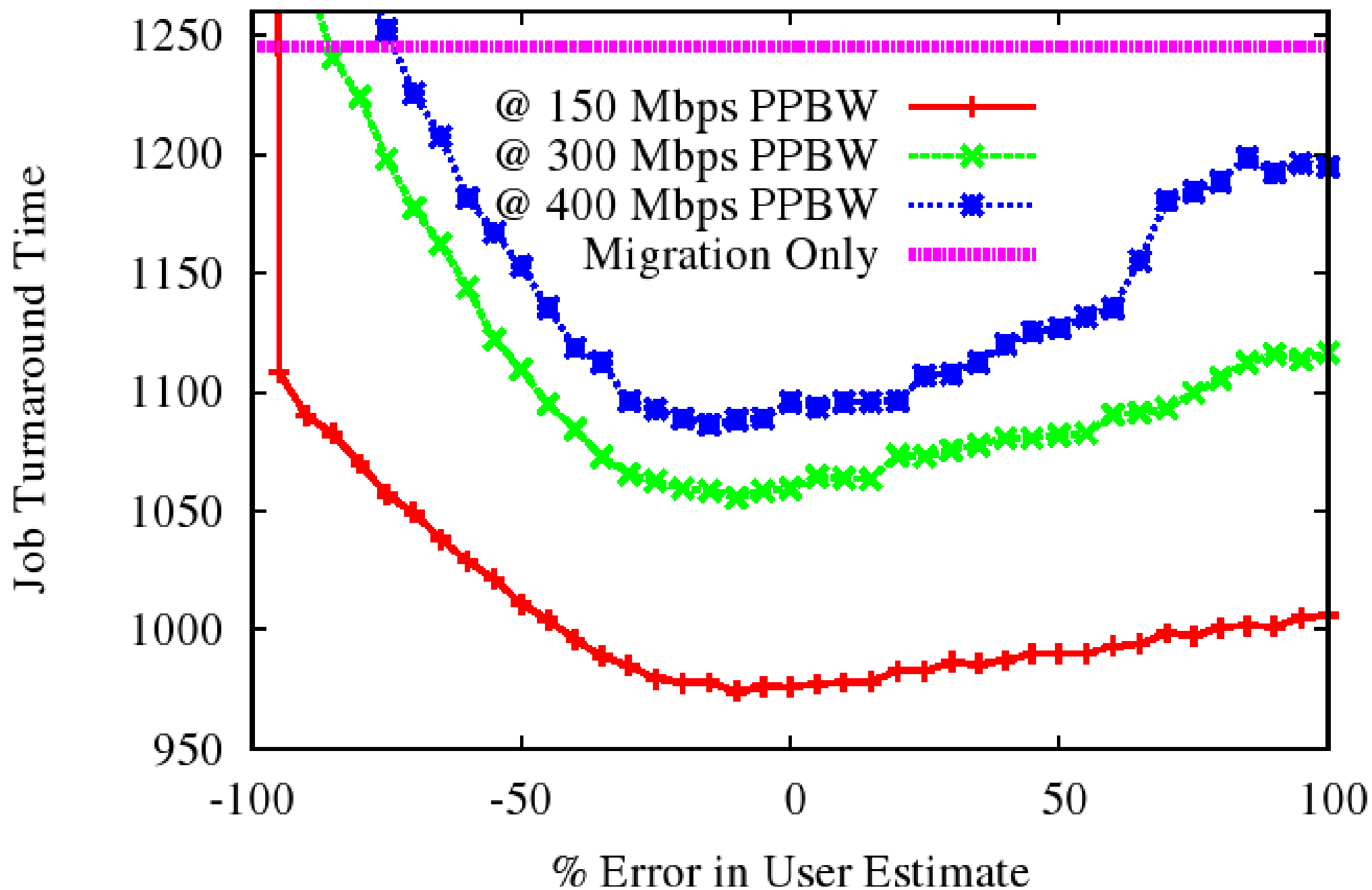
Turnaround Time vs Info vs Error w/ PPBW = 400



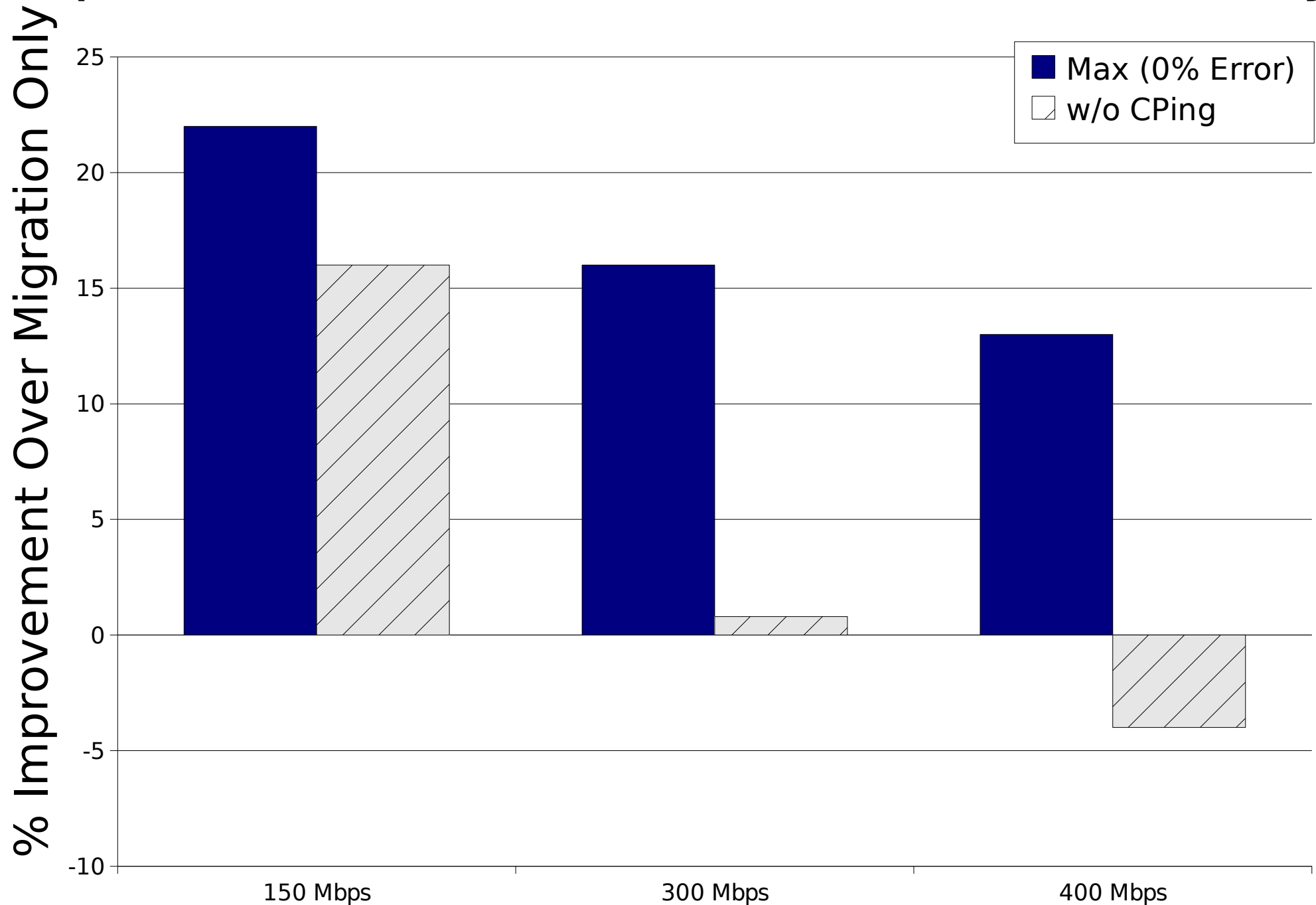
Turnaround Time vs Info w/ PPBW = 400 Mbps



Turnaround Time vs Error w/ Info = 100%



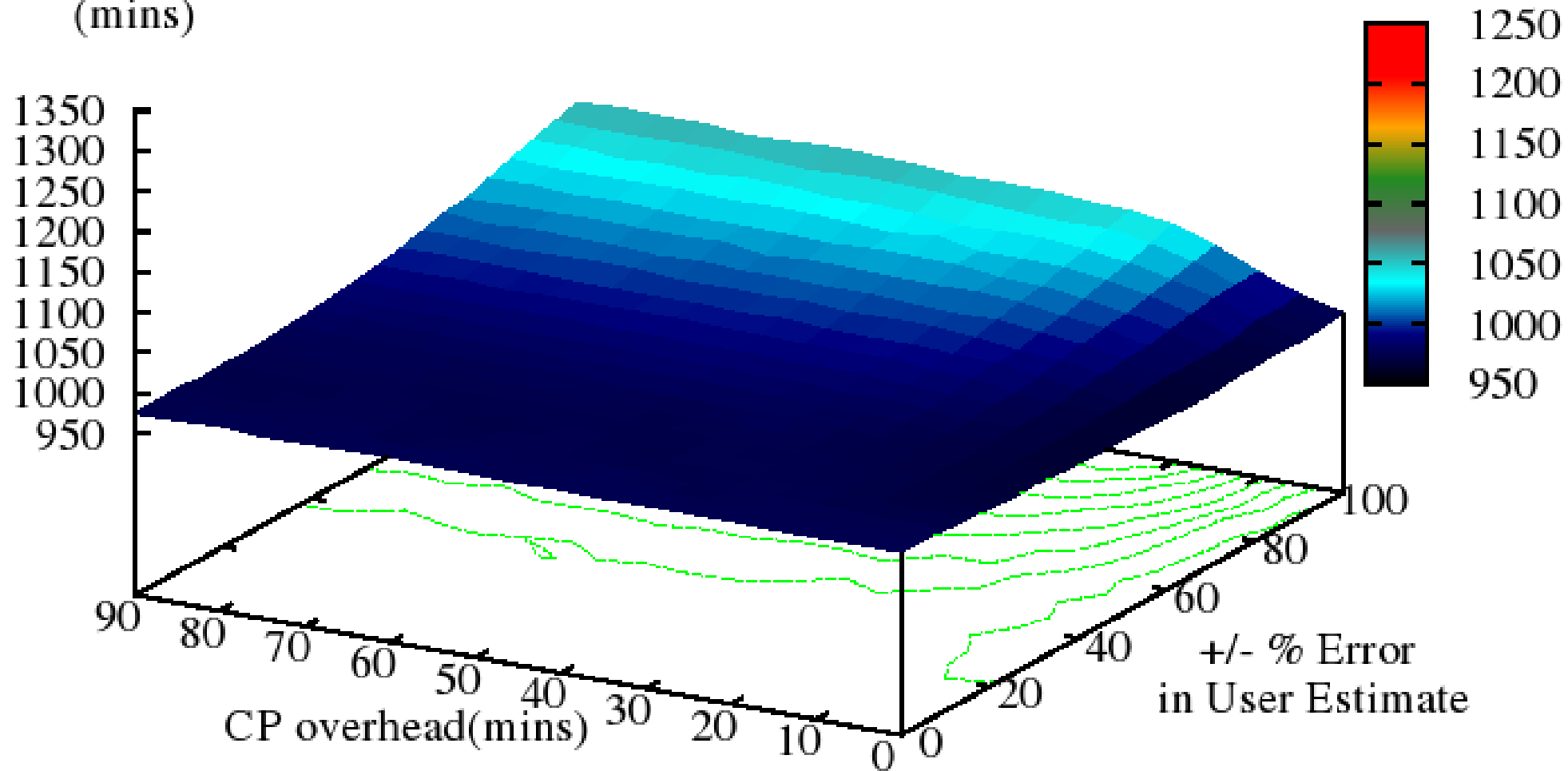
Impact of +/- 100% Estimate Inaccuracy



What potential gain might there be to employ **checkpointing** and **run-time job migration** to mitigate network over-subscription?

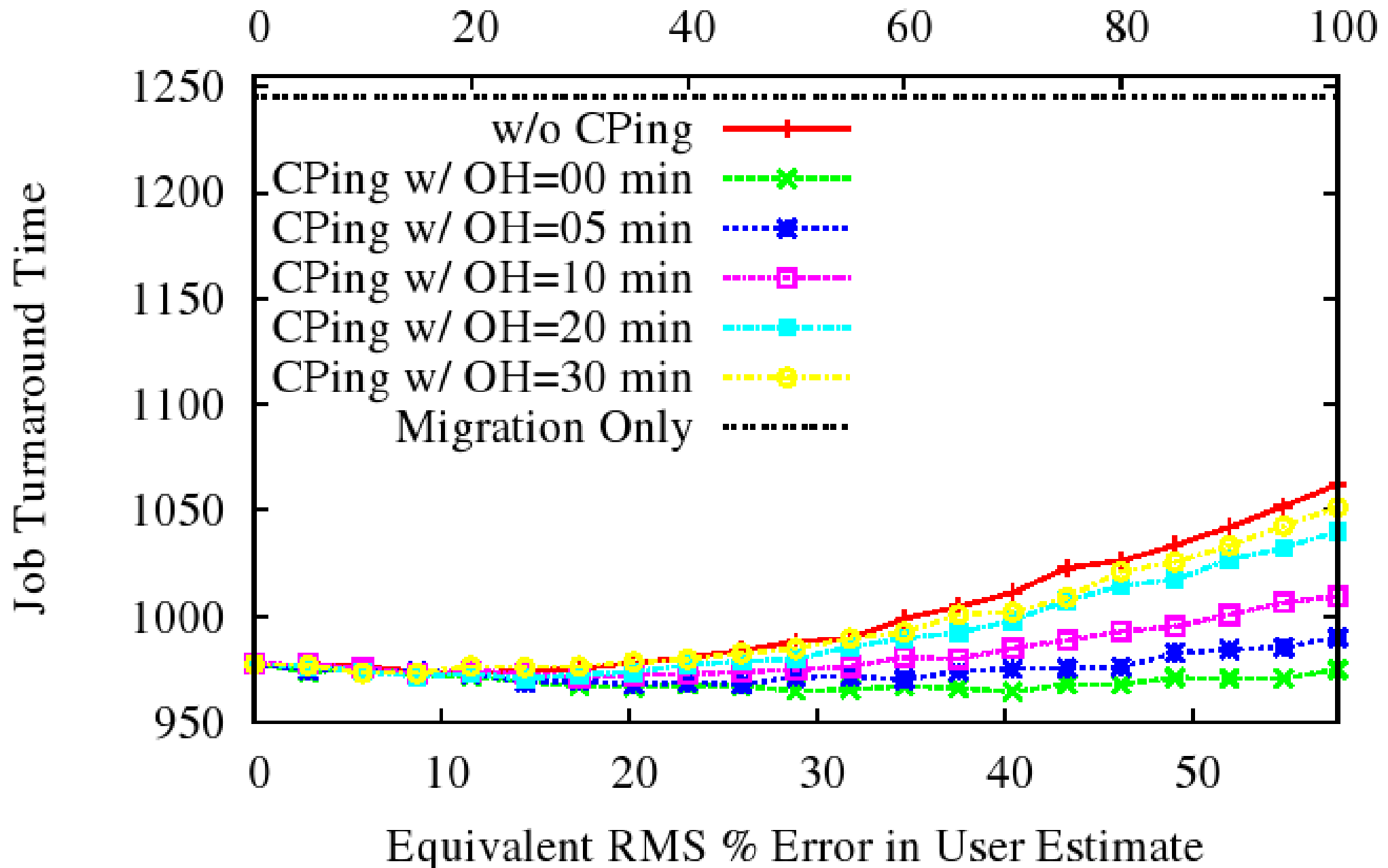
Turnaround Time vs CP Overhead vs Error w/ PPB W = 150 Mbps

Job Turnaround Time
(mins)



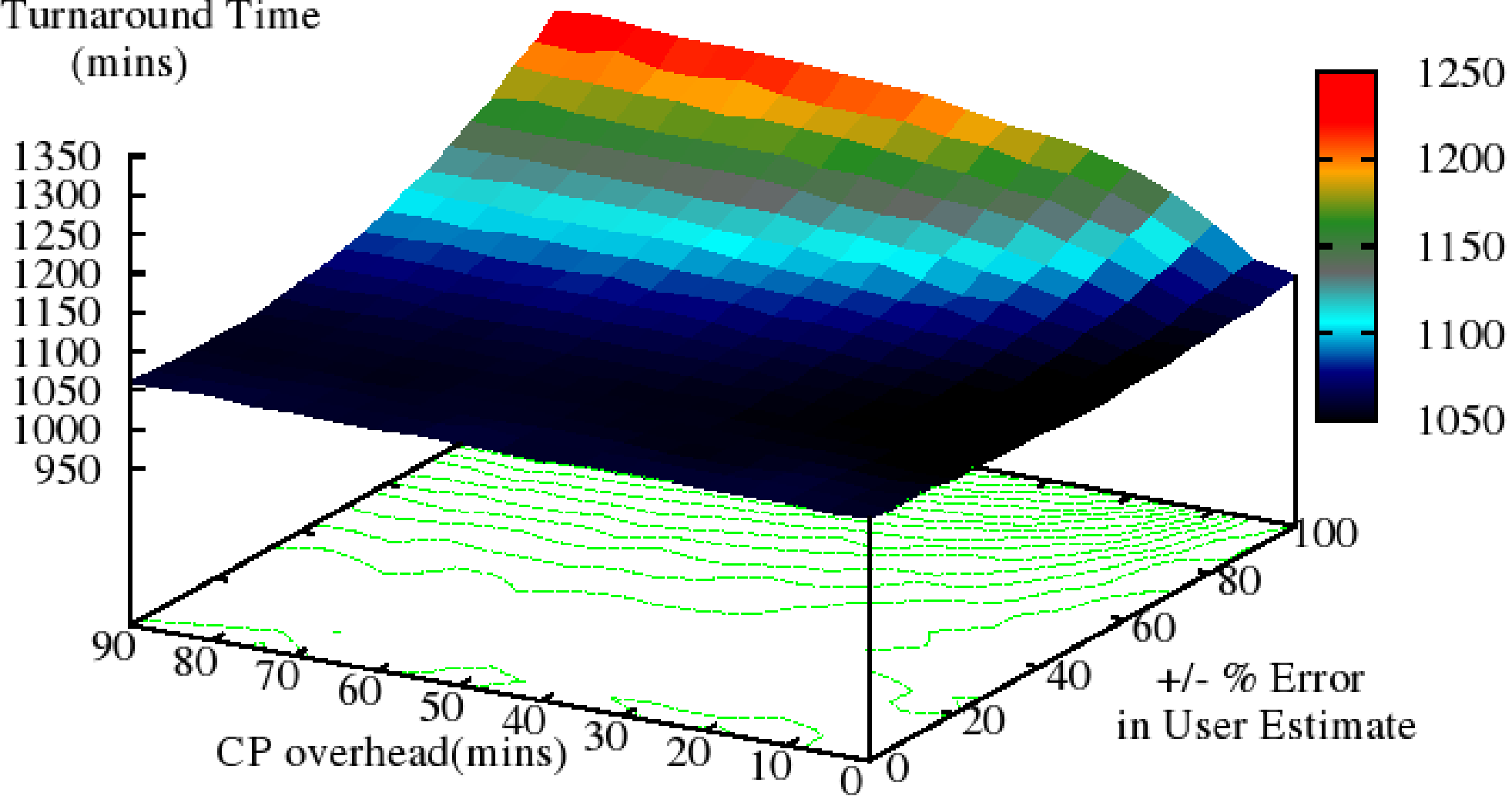
Turnaround Time vs Error w/ PPBW = 150 Mbps

+/- % Error in User Estimate



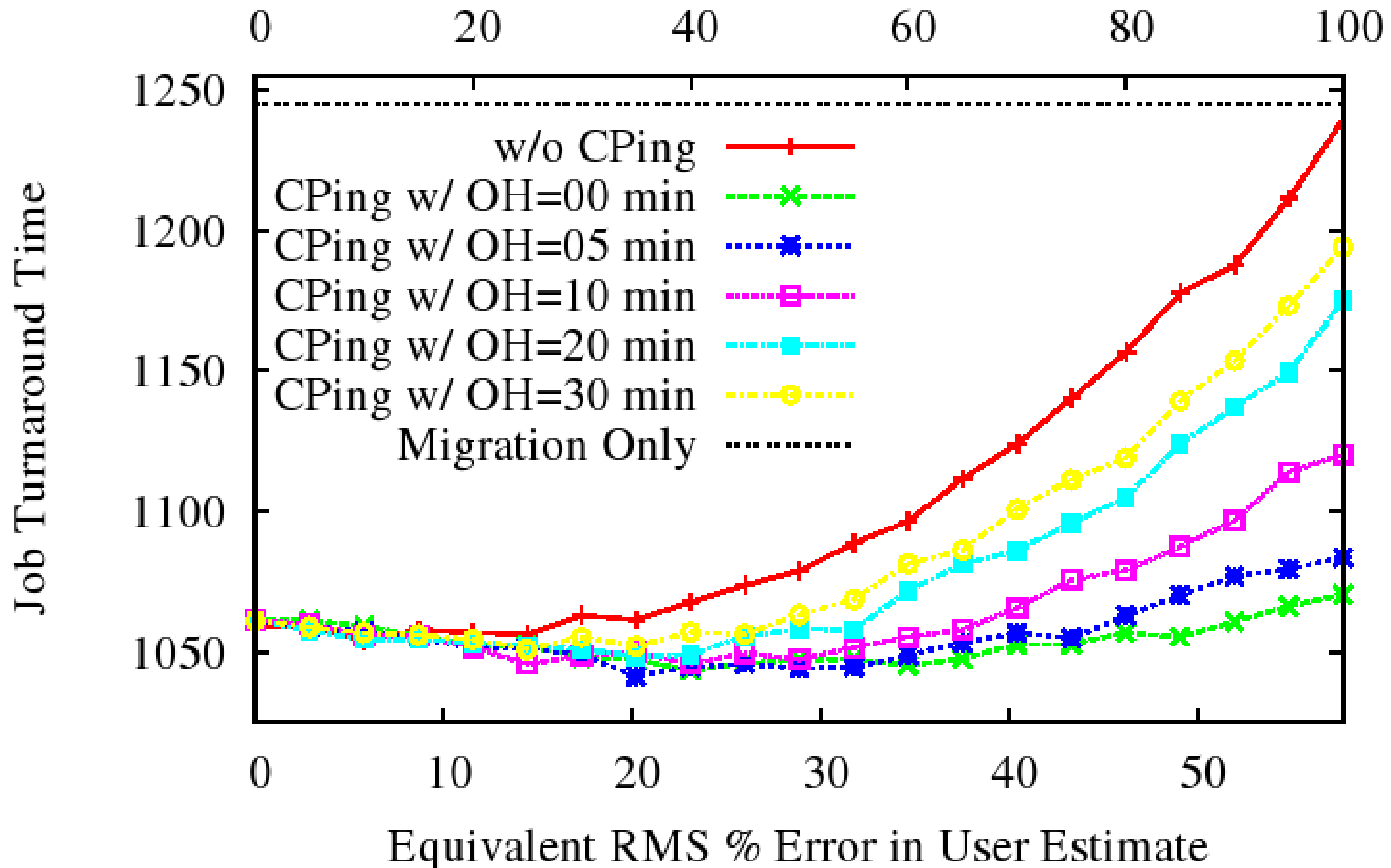
Turnaround Time vs CP Overhead vs Error w/ PPBW = 300 Mbps

Job Turnaround Time
(mins)



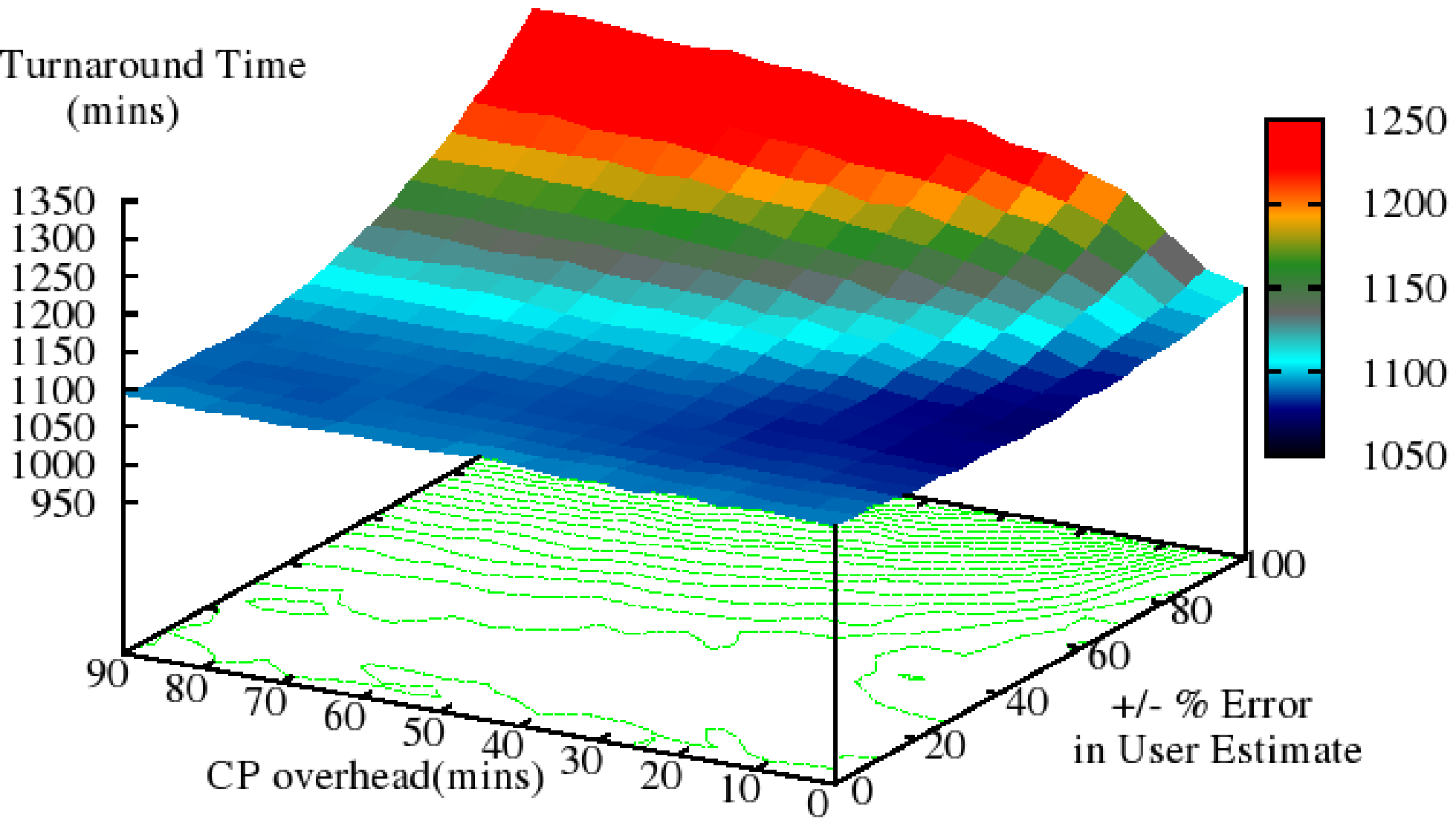
Turnaround Time vs Error w/ PPBW = 300 Mbps

+/- % Error in User Estimate



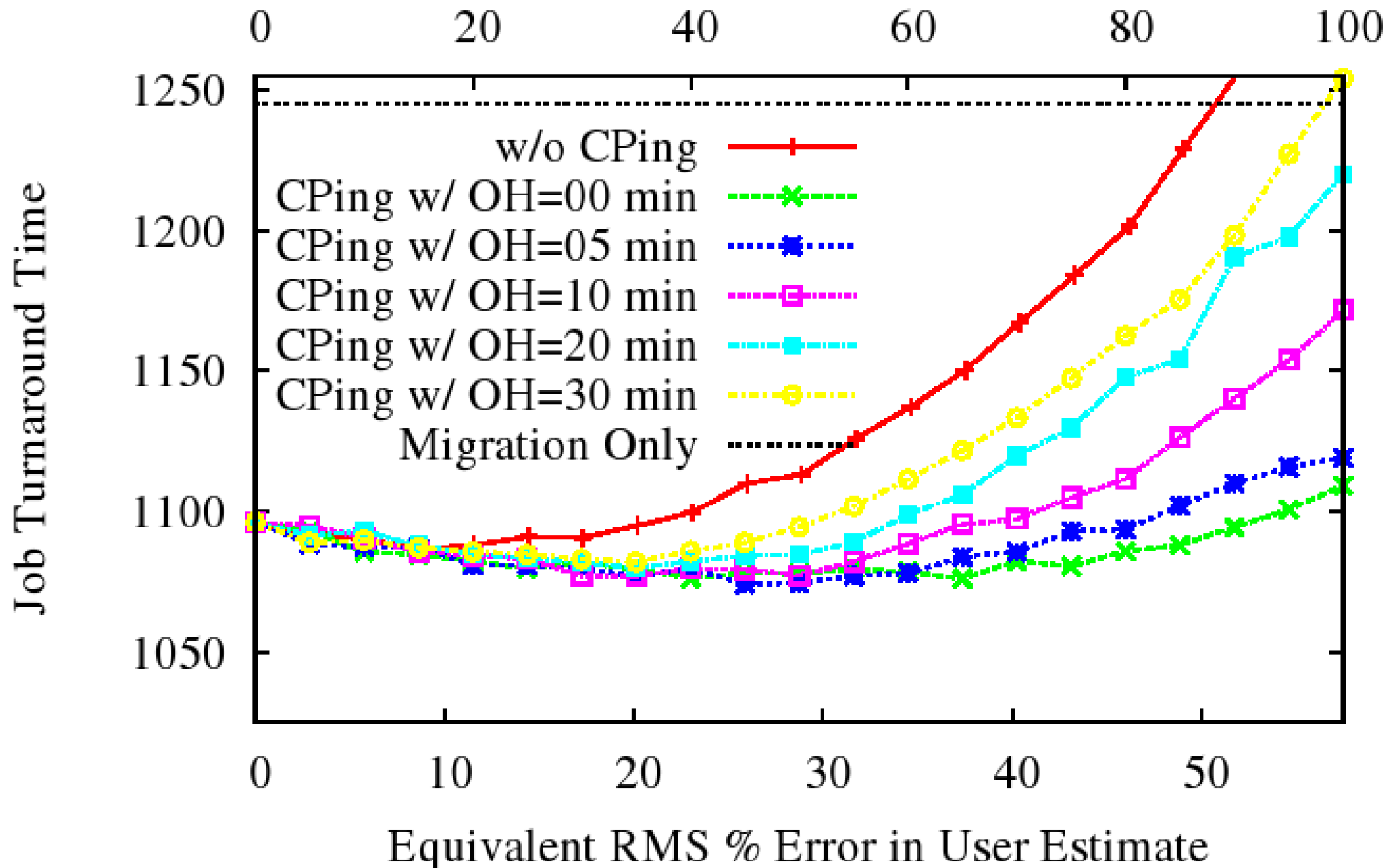
Turnaround Time vs CP Overhead vs Error w/ PPBW = 400 Mbps

Job Turnaround Time
(mins)

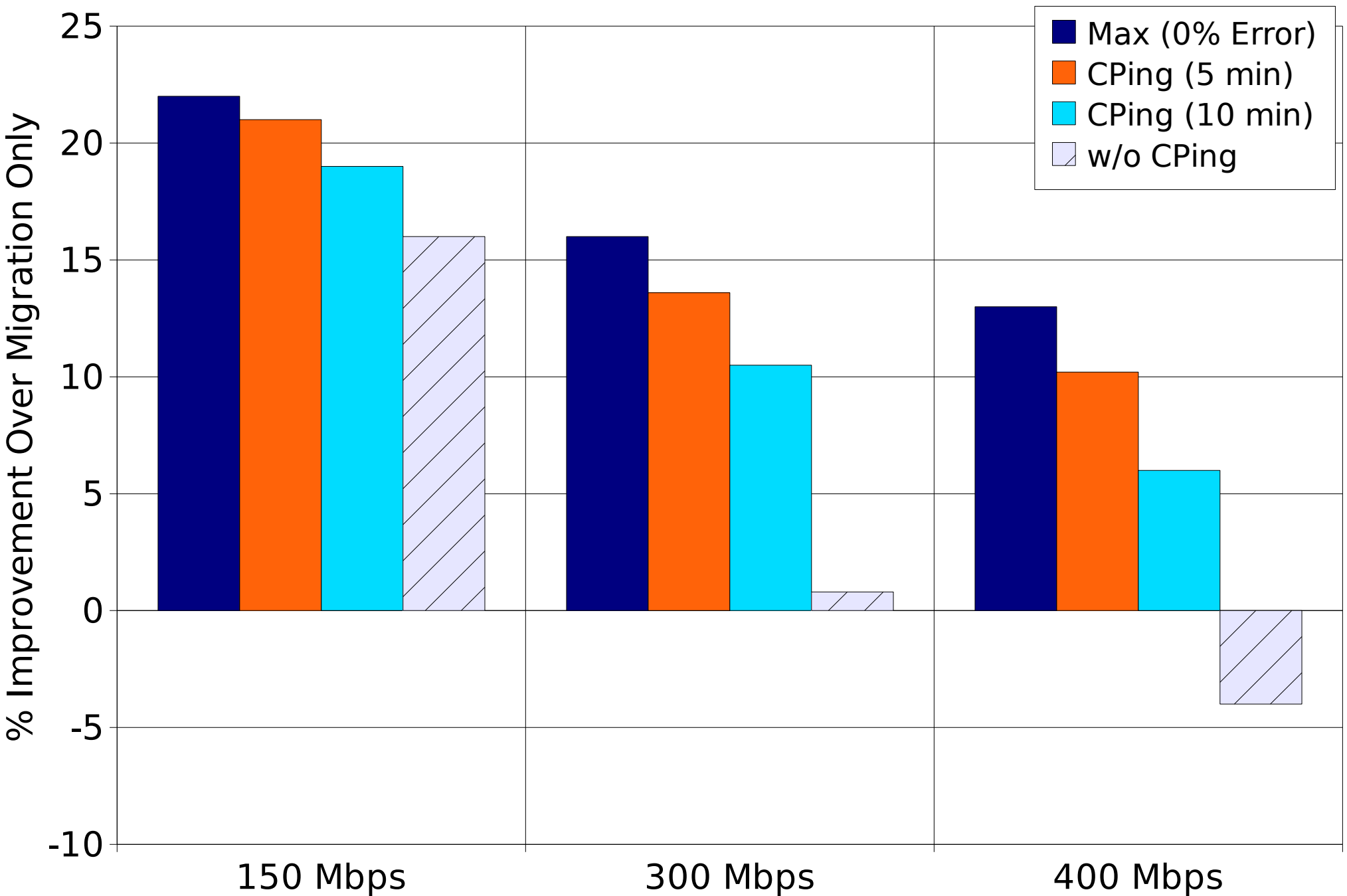


Turnaround Time vs Error w/ PPBW = 400 Mbps

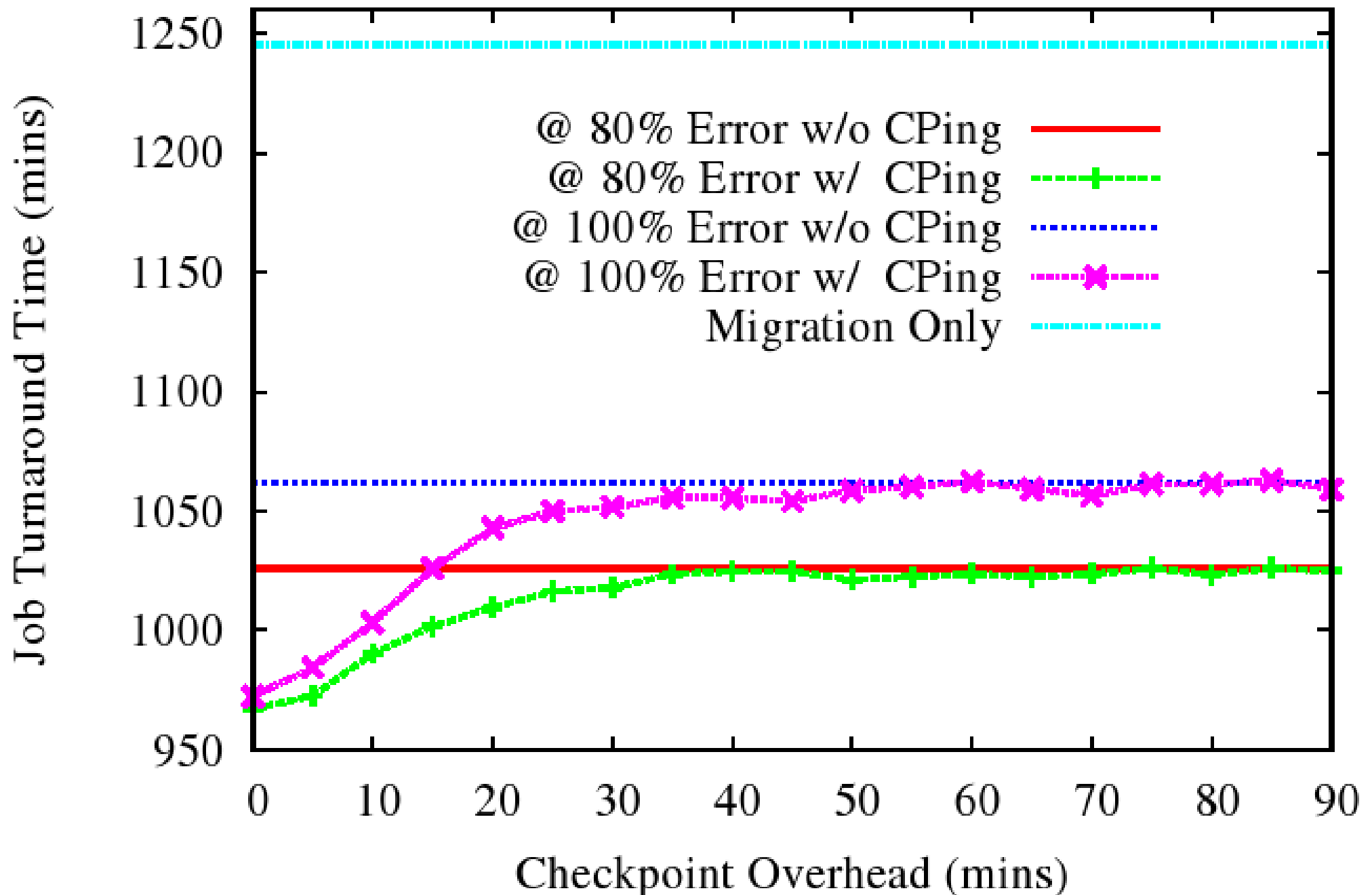
+/- % Error in User Estimate



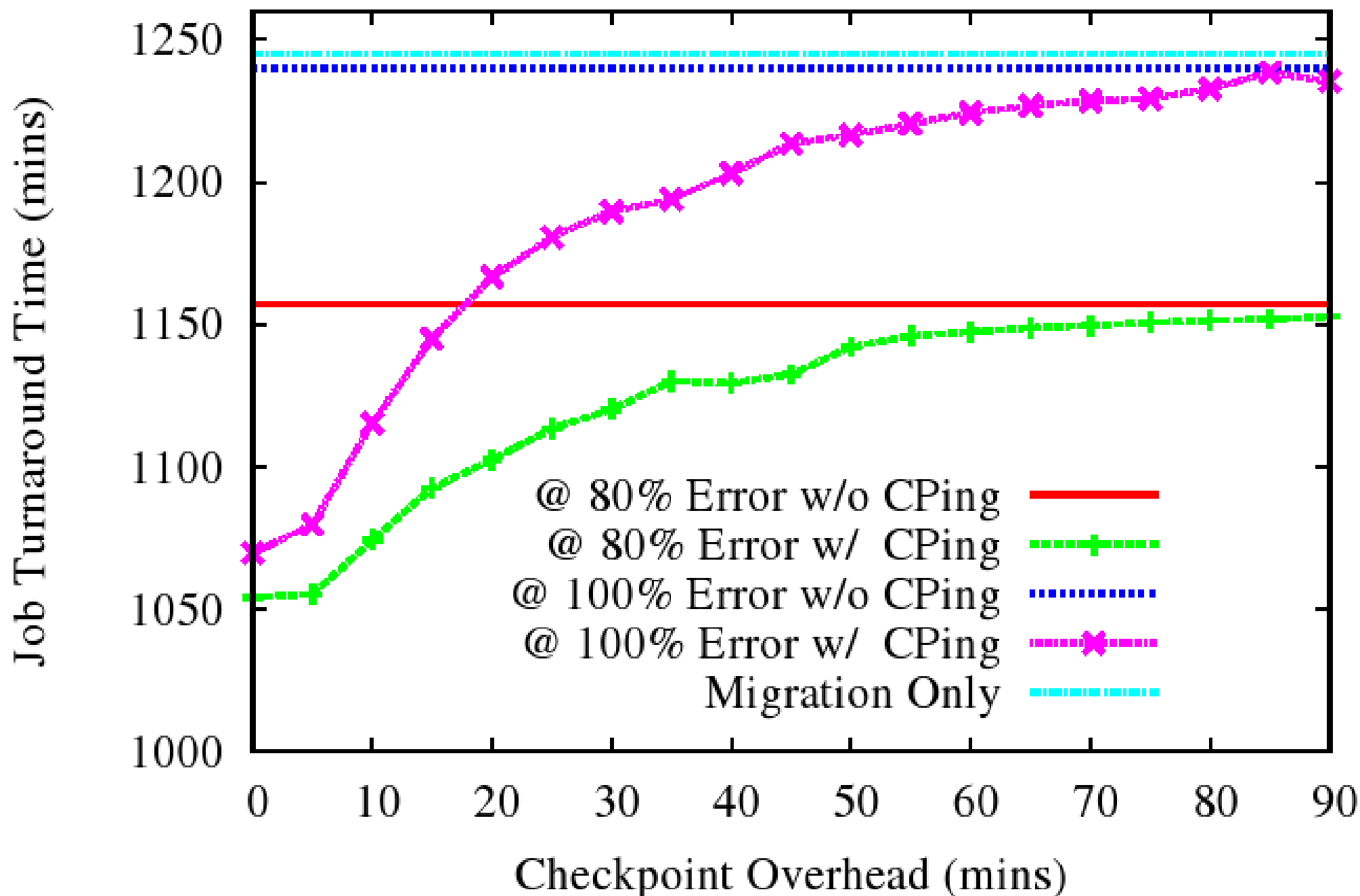
Results Summary at 5 and 10 Minute CP Overhead at 100% Estimate Error



Turnaround Time vs CP Overhead w/ PPBW=150 Mbps



Turnaround Time vs CP Overhead w/ PPBW=300 Mbps





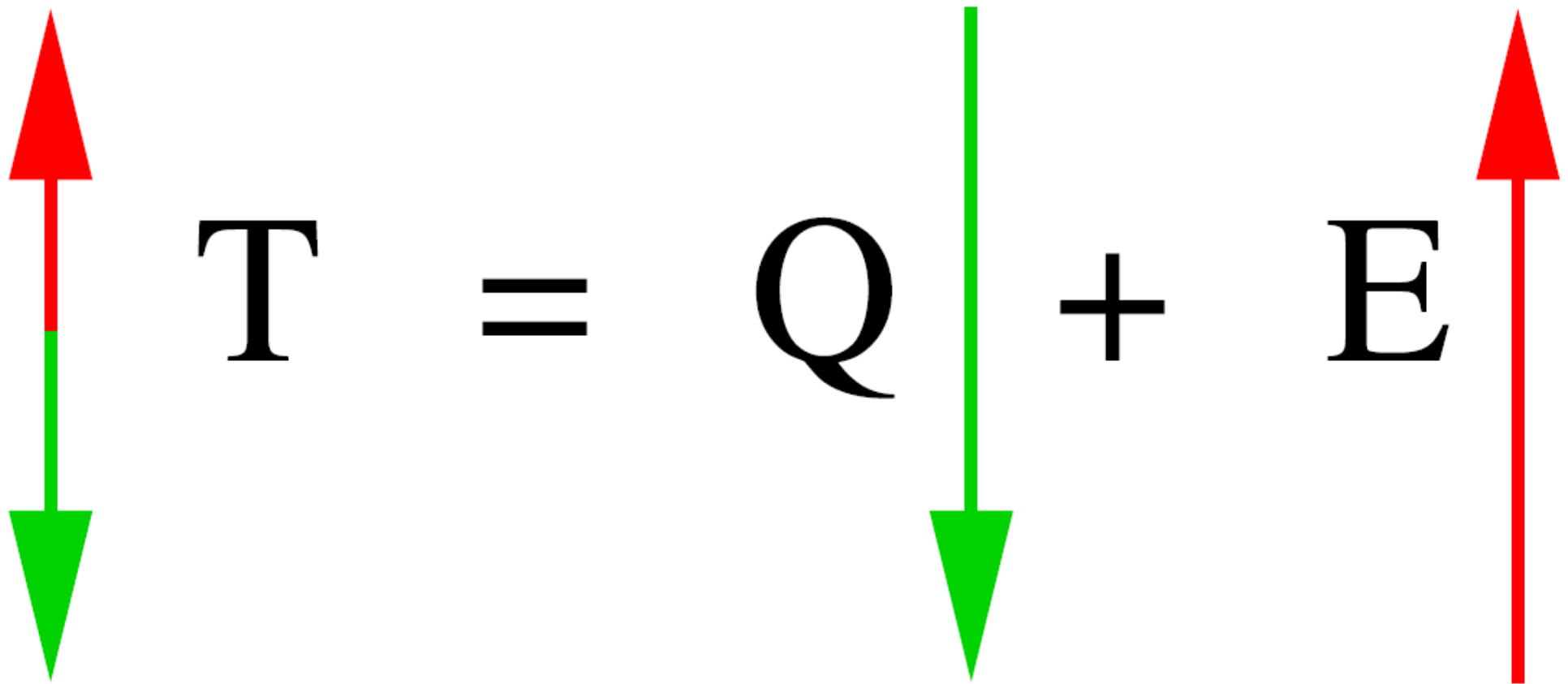
Thank you!

Questions?

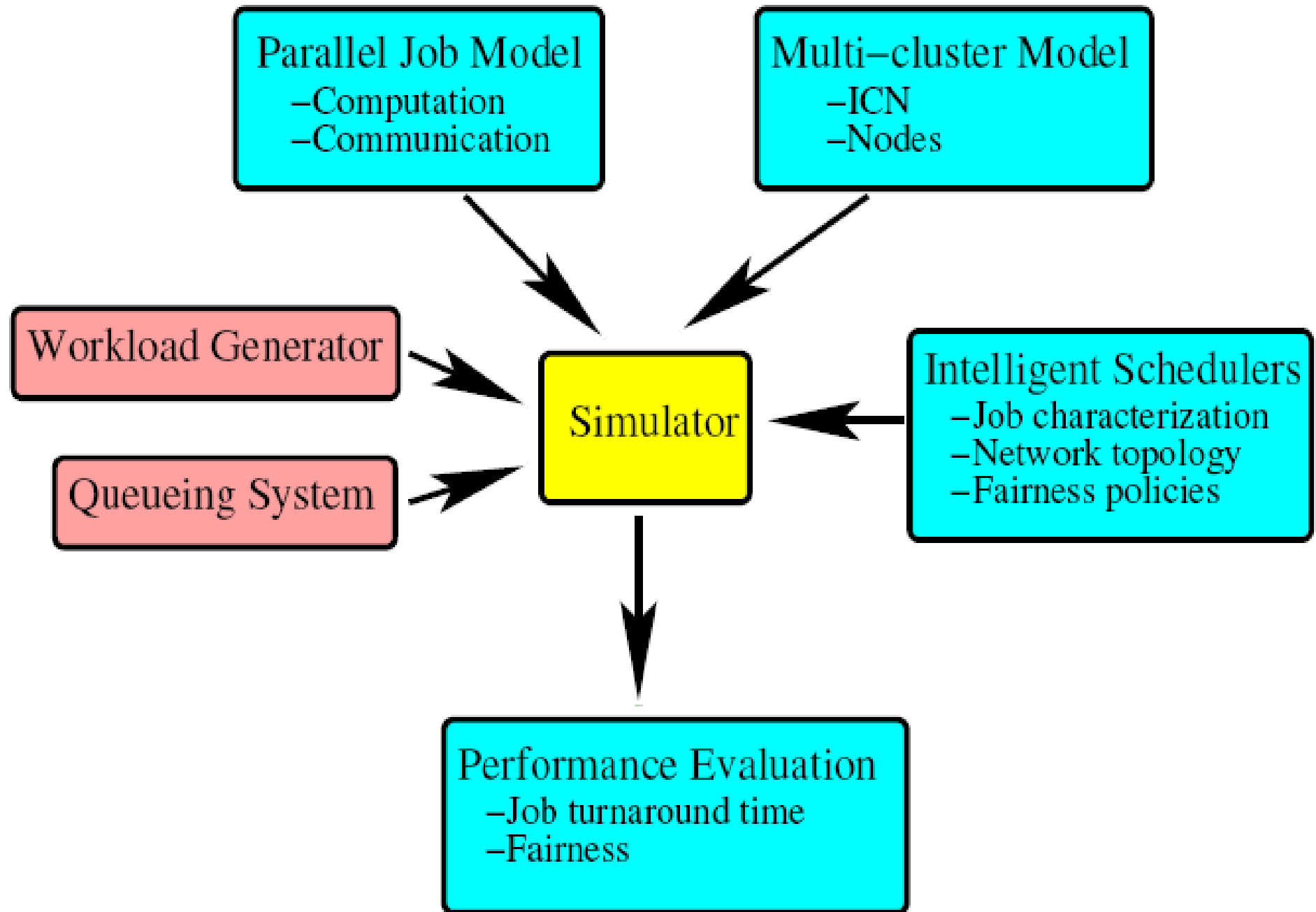
William M. Jones

<http://www.parl.clemson.edu/beosim>

What could some error be 'good'?

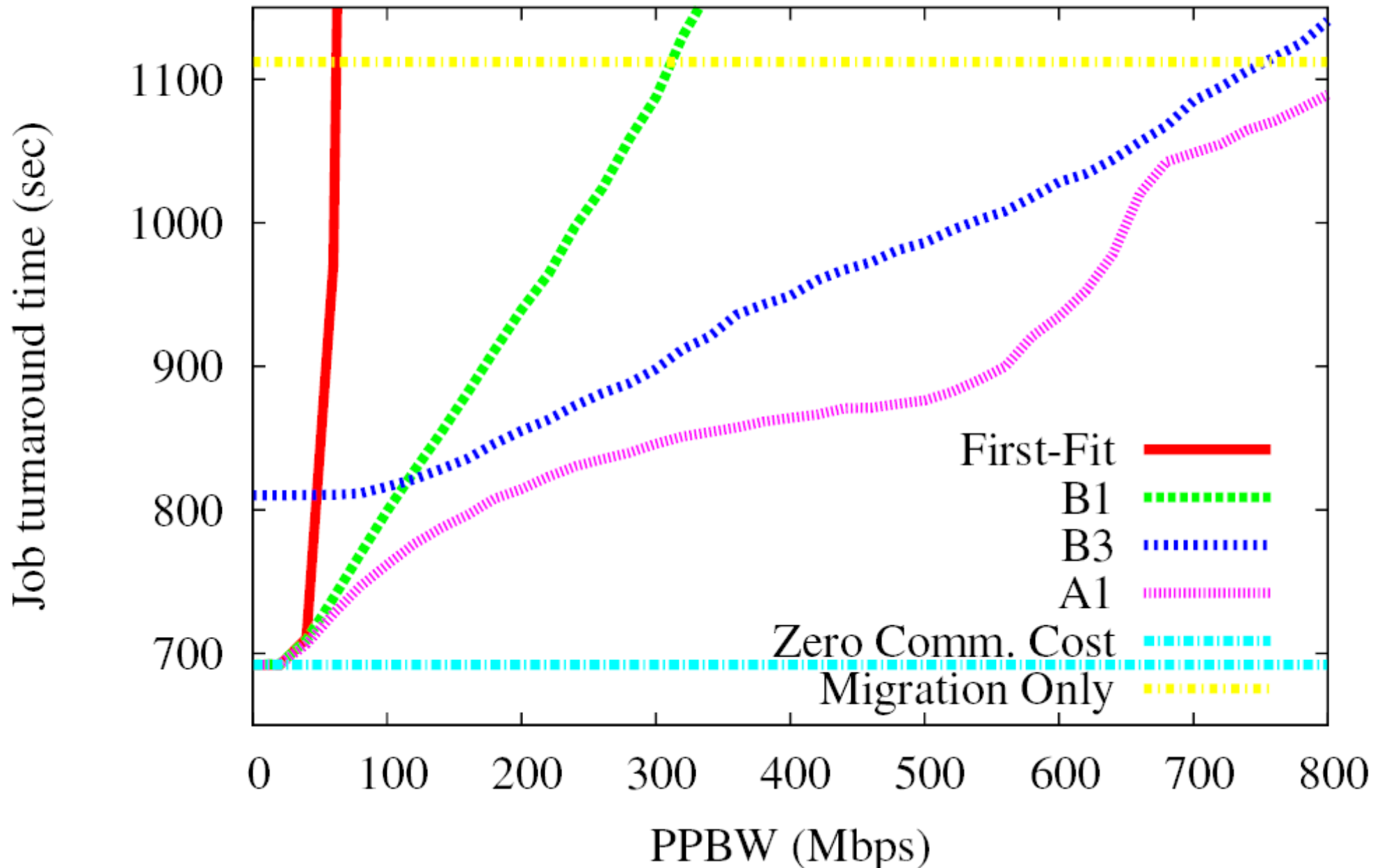


Simulation Framework



Initial Results

Algorithm Comparison, Synthetic Workload

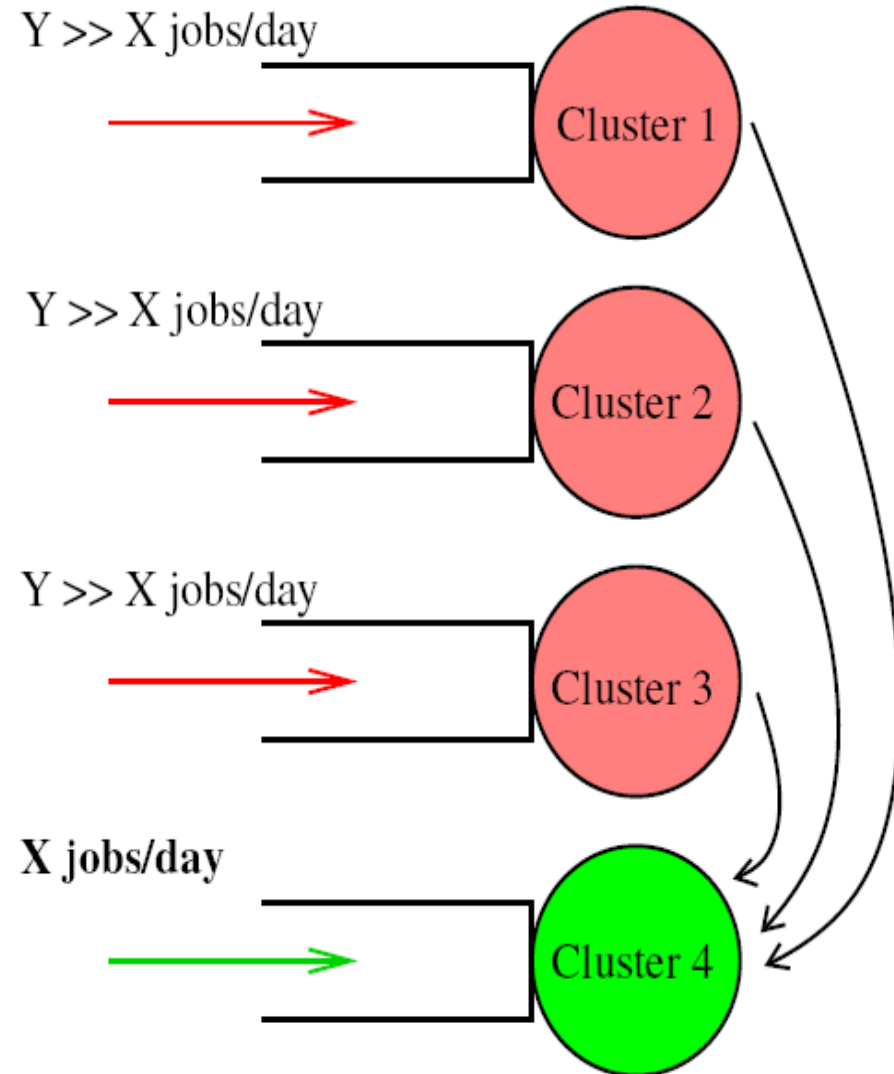


Algorithm Run-Time Analysis

Module	Complexity	Time (μSec)
A1	$O(n \cdot p^m)$	69.6
B1	$O(m \cdot \log(m) + n)$	0.98
B2	$O(m \cdot \log(m) + n)$	0.89
B3	$O(m \cdot \log(m) + n)$	1.28
B4	$O(p)$	1.10

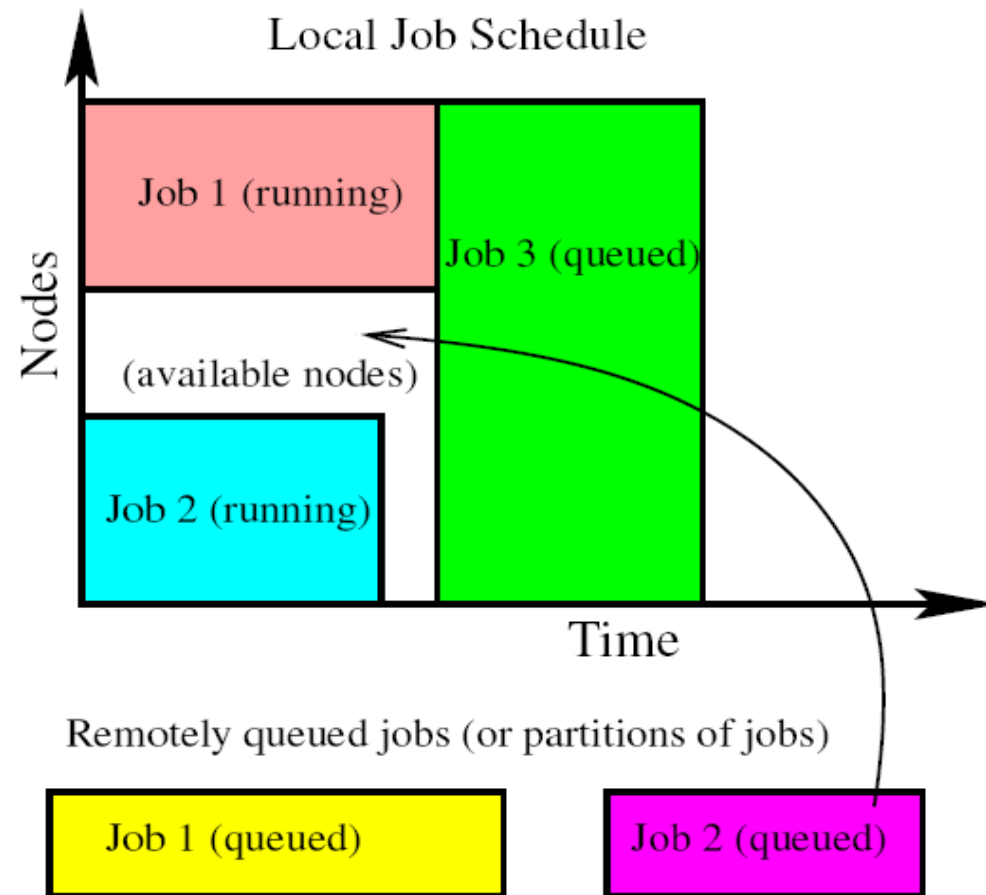
Ensuring Fairness

- Disparate workload intensities
- Different cluster sizes
- Unfair resource sharing
- Overload remote clusters
- Worse than not participating
- Technique to control fairness

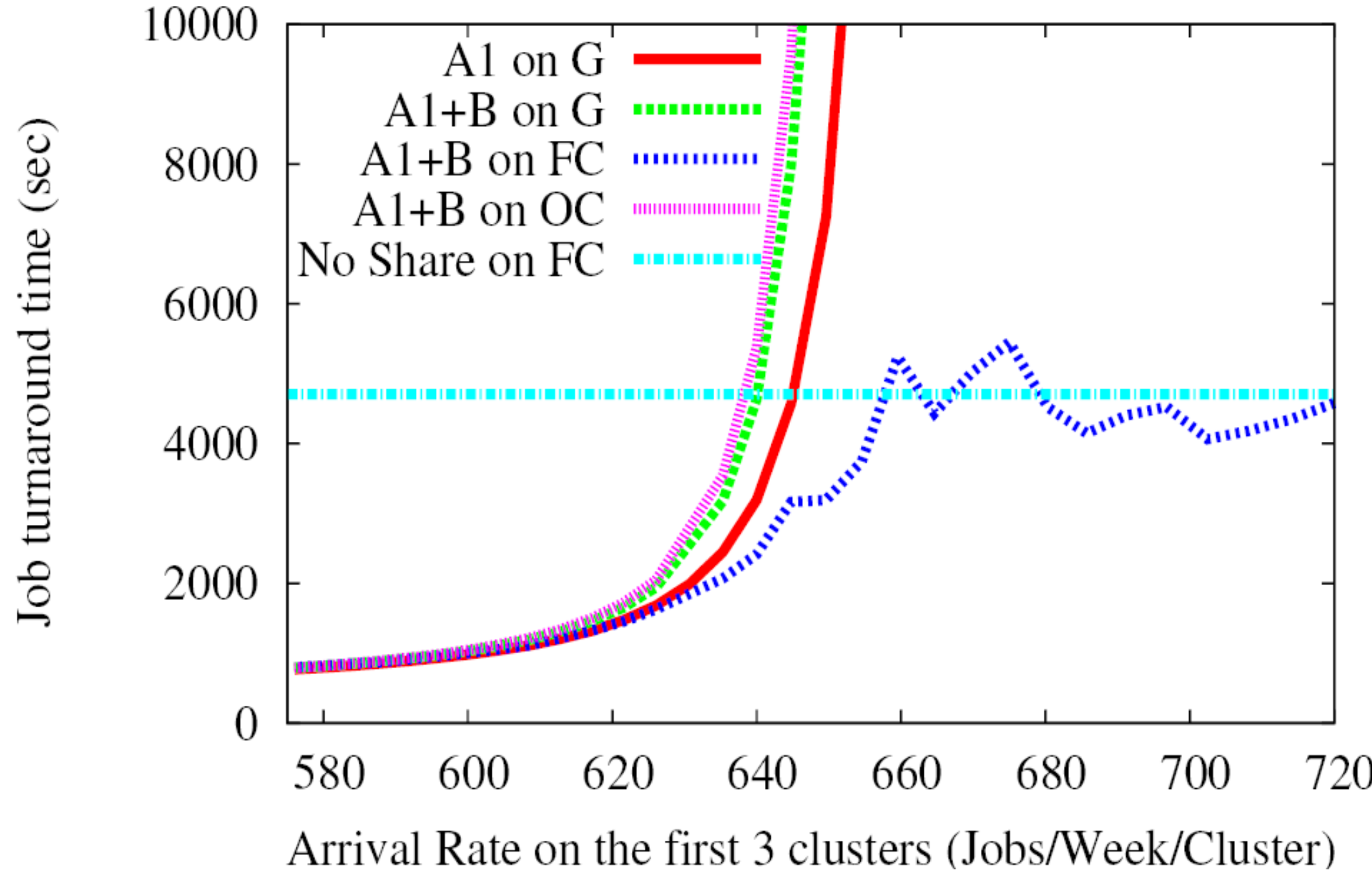


Fairness Via Conservative Backfilling

- Out of order execution
- **No delay to start time**
- Two-tiered approach
- Backfill local jobs first
- Local backfill schedules
- Consider remote job backfill
- Constrain non-local node use
- **Prevents local job starvation**



Fairness, 3 Clusters (100, 100, 100) w/ Increasing Load, 1 Cluster (100) w/ Fixed Load



What about fault-tolerance?

Suppose you could detect that an error occurred, migrate the job, and restart the job from last checkpoint.

How quickly would you need to determine that an interrupt occurred?

Turnaround Time vs CPdelta

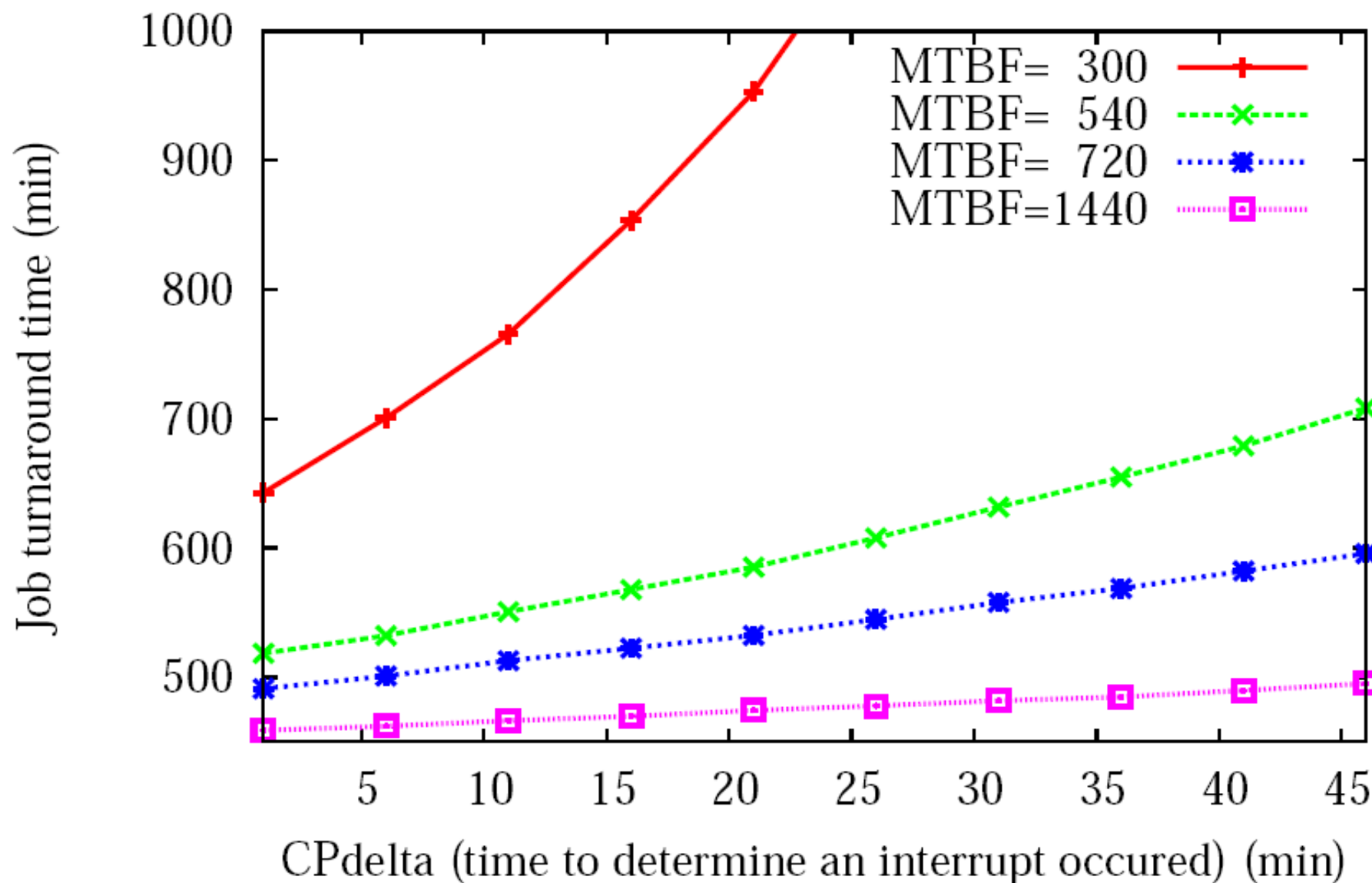


Fig. 4. Turnaround time as a function of CPdelta at various failure rates. *Note that as failures become more frequent, the overall performance becomes more sensitive to reductions in the detection latency, CPdelta.*

Interrupted Job Execution Time vs CPdelta

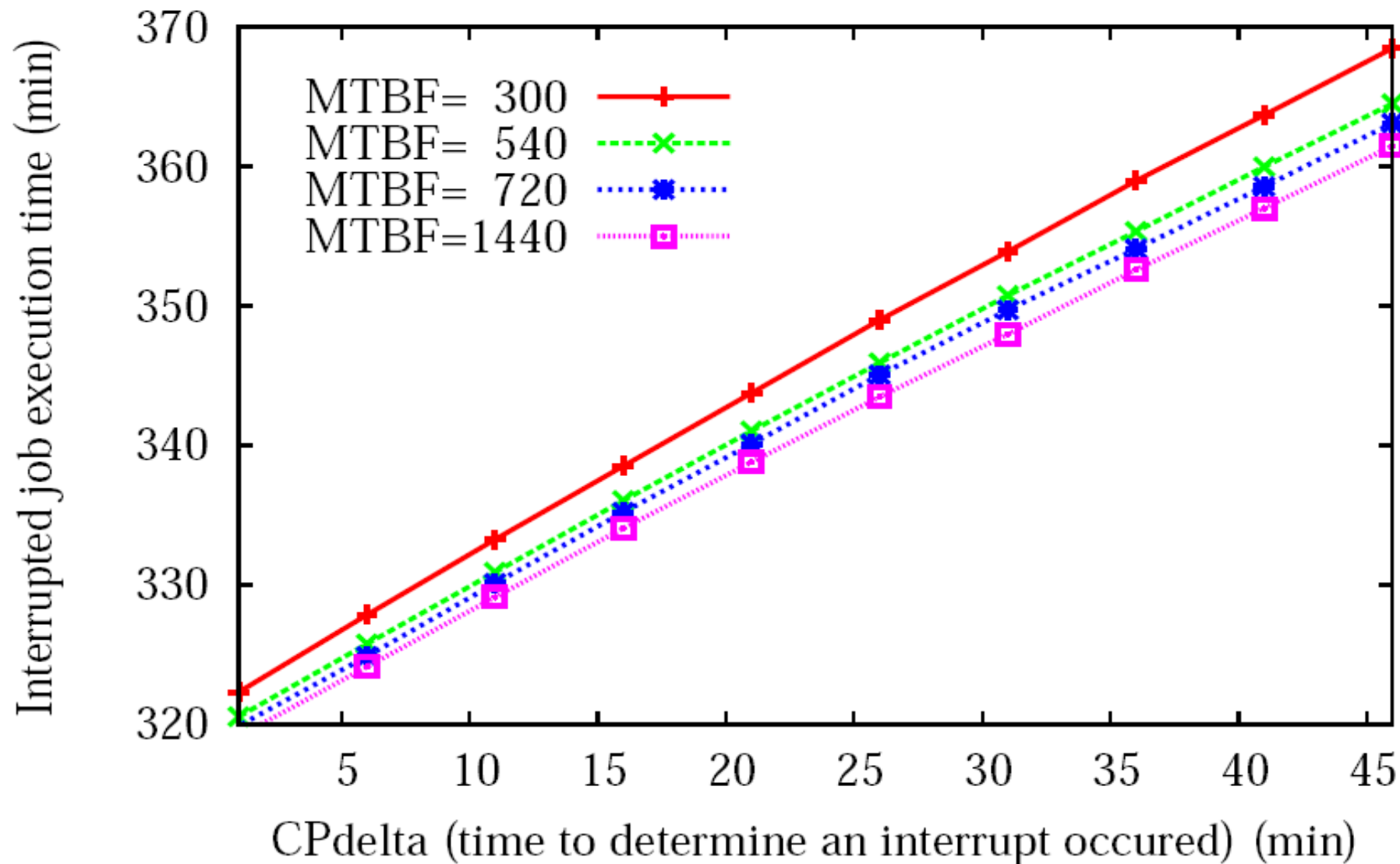


Fig. 6. Job runtime as a function of CPdelta. Here the initial job runtime distribution is the same for all classes of jobs. This is done to illustrate the expected relationship between MTBF and execution time due to multiple single-application interruptions. *Note the similar improvement due to a reduction in CPdelta.*